RESEARCH

Open Access

eSVD-DE: cohort-wide differential expression in single-cell RNA-seq data using exponential-family embeddings



Kevin Z. Lin^{1*}, Yixuan Qiu² and Kathryn Roeder³

*Correspondence: kzlin@uw.edu

¹ Department of Biostatistics, University of Washington, Seattle, WA, USA
² School of Statistics and Management, Shanghai University of Finance and Economics, Shanghai, People's Republic of China
³ Department of Statistics and Data Science, Carnegie Mellon University, Pittsburgh, PA, USA

Abstract

Background: Single-cell RNA-sequencing (scRNA) datasets are becoming increasingly popular in clinical and cohort studies, but there is a lack of methods to investigate differentially expressed (DE) genes among such datasets with numerous individuals. While numerous methods exist to find DE genes for scRNA data from limited individuals, differential-expression testing for large cohorts of case and control individuals using scRNA data poses unique challenges due to substantial effects of human variation, i.e., individual-level confounding covariates that are difficult to account for in the presence of sparsely-observed genes.

Results: We develop the eSVD-DE, a matrix factorization that pools information across genes and removes confounding covariate effects, followed by a novel two-sample test in mean expression between case and control individuals. In general, differential testing after dimension reduction yields an inflation of Type-1 errors. However, we overcome this by testing for differences between the case and control individuals' posterior mean distributions via a hierarchical model. In previously published datasets of various biological systems, eSVD-DE has more accuracy and power compared to other DE methods typically repurposed for analyzing cohort-wide differential expression.

Conclusions: eSVD-DE proposes a novel and powerful way to test for DE genes among cohorts after performing a dimension reduction. Accurate identification of differential expression on the individual level, instead of the cell level, is important for linking scRNA-seq studies to our understanding of the human population.

Keywords: Case–control subjects, Gamma–Poisson distribution, Matrix factorization, Multi-individual data

Background

High-throughput single-cell RNA-seq (scRNA) technology has advanced tremendously over the last decade and helped biologists uncover differing cell-type proportions as well as differentially expressed (DE) genes within a particular cell-type when studying various diseases or disorders. These findings were previously inaccessible using bulk RNA-seq technology. As the technology has developed more in accuracy and cost-efficiency, many



© The Author(s) 2024. **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/publicdommain/zero/1.0/) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

labs have begun sequencing entire cohorts of individuals to study up-regulated/downregulated pathways specific to each cell type in diseases/disorders that generalize to the human population. For example, biologists have sequenced hundreds of thousands of cells for 44 individuals with lung adenocarcinoma [1], hundreds of thousands of neurons for 84 individuals with varying severity of Alzheimer's Disease [2], and millions of blood cells for 162 individuals with systemic lupus erythematosus and 99 control individuals [3] in order to discover DE genes for the disease/disorder consistent with the entire cohort. While numerous existing and benchmarked single-cell DE methods are designed to find differentially-expressed patterns among cells [4, 5], this task fundamentally differs from the cohort-wide studies' goal of finding differentially expressed patterns among *individuals*. This pressing methodological gap has been raised by many biologists who have collected cohort-wide scRNA-seq datasets [6, 7]. Hence, we reinvestigate the shortcomings of existing DE methods commonly used to analyze cohort-wide scRNAseq data and design a new DE method specifically suited for such data. Focusing on scRNA-seq data enables us to study more principled ways to model cohort-wide scRNAseq data, which we hope could be used to inspire cohort-wide DE methods for other single-cell technologies in the future, as well as pioneer applications in future eQTL, as well as upcoming studies that sequence spatial transcriptomics or paired multiomics on large cohorts.

One of the main distinctions between a DE analysis among cells compared to a DE analysis among individuals lies in how the case and control population distributions are quantified. DE analyses among individuals need to account for variability within and among individuals in order to properly model human variation. However, most existing DE methods lack one of these two aspects. On the one hand, variability within an individual hinders "pseudobulk" analyses where all the cells among each individual are summed to yield a pseudobulk sample. Then, methods originally designed for bulk RNA-seq data such as DESeq2 [8] and edgeR [9] are used, but these methods do not account for the variability within an individual. On the other hand, variability across individuals hinders most DE methods made for scRNA-seq data. Specifically, a gene could be differentially expressed among cells but not among individuals if all cells with a significantly higher gene expression come from a small subset of individuals.

The second main distinction is that in cohort-wide scRNA-seq data, there are potentially substantial effects of individual-level covariates such as age, sex, and smoking status that can induce differences in gene expression among the cells that do not reflect the biological differences related to the disease or disorder. A conventional strategy is to regress out the covariate effects for each gene one at a time via a mixed effects model such as MAST [10] and NEBULA [11] where there is a random effect for each individual. However, this regression might be inaccurate since each gene is sparsely sequenced, detrimentally impacting the downstream DE analysis. Hence, an alternative strategy is to use a dimension-reduction method via matrix factorization such as GLM-PCA [12], ZINB-WaVE [13], or scGBM [14]. These methods pool information across genes to remove confounding covariates' effects more effectively. However, a naive application of DE testing on the dimension-reduced scRNA-seq data has been observed to inflate the Type-1 error [15]. This inflation occurs because the dimension reduction introduces correlations among genes that contaminate the signal. Null genes could seem significantly (and erroneously) differentially expressed after a dimension reduction if the test is not performed with care. (See an illustration of this phenomenon in Additional file 1: Fig. A1.) The main focus of our paper is to solve this statistical dilemma of how to perform differential testing after dimension reduction. We note that another recent line of work has developed differential testing among cells after a dimension reduction via deep variational autoencoders [16, 17]. The authors overcome the aforementioned pitfall of DE testing after dimension reduction by leveraging the inherent randomness in autoencoders and adding pseudocounts to avoid deeming genes with a small log-fold change as significant. However, we do not pursue this direction since matrix factorizations offer practitioners a more transparent and interpretable framework.

We design the exponential-family SVD differential expression (eSVD-DE) to overcome these two main obstacles, which extends our previous work [18]. Importantly, our method infers the differential expression based on the posterior distribution after performing a dimension reduction, which helps counteract the Type-1 error inflation. This combination of matrix factorization, the posterior distribution, and a test statistic designed to assess differential expression among individuals enables eSVD-DE to better detect DE genes in cohort studies compared to current methods. The eSVD-DE also enables model diagnostics to assess if the assumed statistical model is appropriate for modeling the scRNA-seq dataset.

In this paper, we focus on testing for cells of a particular cell type. We show that eSVD-DE can find reproducible signals in multiple pairs of cohort datasets, either across various cell types between two independent studies of idiopathic pulmonary fibrosis (IPF) in the human lung [6, 19] or within a study of non-inflamed and inflamed cells studying ulcerative colitis in the human colon [20]. We also show that eSVD-DE can find novel DE genes across different cell types in a dataset studying autism [21]. Altogether, these analyses provide evidence that eSVD-DE is a valuable tool for investigating differential expression among cohort studies that will become more prevalent as high-throughput sequencing technologies are applied to large cohorts.

Results

Overview of eSVD-DE for differential expression testing

eSVD-DE performs DE by first projecting the cells onto a low-dimensional manifold while removing the effects of covariates. This step is the cornerstone and namesake of our method. Our dimension reduction follows previous dimension-reduction work such as GLM-PCA [12], ZINB-WaVE [13], and scGBM [14], where we embed the cells via the Poisson distribution based on the gene expression $A \in \{0, 1, 2, ...\}^{p \times n}$ and the covariate matrix $C \in \mathbb{R}^{n \times r}$, where n, p, and r denote the number of cells, genes, and covariates (Fig. 1A). Importantly, these covariates contain an intercept term, the log sequencing depth (computed as the log of the total counts per cell), the case–control indicator, and covariates that could be potential confounders, such as clinical covariates of each individual (sex, age, and smoking status). We denote the specific covariate for the case–control indicator as $C_{\cdot,(cc)}$. The eSVD-DE learns a coefficient matrix that removes the effects of the covariates as well as the low-dimensional embedding of "residuals" X via the hierarchical model for gene $j \in \{1, ..., p\}$ and cell $i \in \{1, ..., n\}$ following work such as [22],



A) Step 1: Low-dimensional embedding via Poisson distribution (using eSVD)

Step 3: Differential expression after removing effects of confounding covariates



Fig. 1 A Schematic of the eSVD-DE's matrix factorization, where the observed scRNA-seq data is modeled as a sum of two low-rank matrices, one for the covariates and one for the cells' latent vectors, with an exponential link function (for the Poisson distribution). B The cells' latent vectors can be used for diagnostic checks, such as visualization via Isomap. C To account for overdispersion and over-fitting of the dimension reduction, shrink each cell via the negative binomial distribution's posterior mean. D Represent each individual by a Gaussian distribution among the individual's cells. E Compute a test statistic analogous to the T-test after aggregating cells from the cases or control individuals in the cohort. C through E are performed for each gene. F Volcano plot, showing a multiple testing cutoff to determine the significant DE genes

 $(A_{ji}|\lambda_{ji}) \sim \text{Poisson}(\ell_{ji} \cdot \lambda_{ji}), \text{ and } \lambda_{ji} \sim \text{Gamma}(\text{mean} = \mu_{ji}; \text{variance} = \gamma_j \cdot \mu_{ji}),$

and

$$\mu_{ji} = \exp\left(\left(Y_{j,\cdot}\right)^{\perp}\left(X_{i,\cdot}\right) + Z_{j,(\mathrm{cc})} \cdot C_{i,(\mathrm{cc})}\right),$$

where $X \in \mathbb{R}^{n \times k}$ and $Y \in \mathbb{R}^{p \times k}$, and ℓ_{ji} denotes the covariate-adjusted sequencing depth that accounts for the effect of all the remaining covariates in *C* aside from $C_{,(cc)}$. Here, γ_j denotes the overdispersion parameter for gene *j*, which measures how much extra variability exists between the assumed Poisson fit in the low-dimensional embedding and the observed data, analogous to works like SAVER [23] and totalVI [24]. While this *k*-dimensional embedding *X* can be helpful for a wide variety of applications downstream, we focus on modifying our previous work [18] to perform DE analysis for cohort data. We advocate using a matrix-factorization-based approach to perform this dimension reduction over deep learning approaches since our approach enables conventional model diagnostics to assess and tune eSVD-DE (Fig. 1B).

After fitting a low-dimensional embedding, we cater the following procedure to test for DE among individuals, an aspect absent from our previous work [18]. First, we adjust the predicted relative expression for each cell's gene expression using its posterior according to the Gamma-Poisson distribution (i.e., Negative Binomial, Fig. 1C). This posterior, denoted by $\lambda_{ji} \mid A_{ji}$, is computed by first estimating the overdispersion γ_j of each gene. Importantly, this posterior distribution is computed to be the relative expression after accounting for the contributions of the confounding covariates. For a particular gene, we summarize the cells' posterior distribution of $\lambda_{ji} \mid A_{ji}$ from each individual as a Gaussian distribution following the Central Limit Theorem, and we then compute the T-test statistic reflecting the collective difference between the Gaussians from case individuals to those from control individuals (Fig. 1D, E). Notably, this means we do not perform a differential expression test based on μ_{ii} 's because different genes are highly correlated based on their values in μ_{ii} due to the low dimensional embedding, which will artificially inflate the Type-1 error. (See the Appendix for a more in-depth discussion.) Finally, we use a multiple-testing procedure based on empirical null distribution to report the DE genes, which has been successful in other settings to account for possible model misspecification [25] (Fig. 1F). More details of the eSVD-DE procedure are described in "Statistical model and method."

eSVD-DE relies on three primary statistical assumptions: (1) scRNA-seq data can be appropriately modeled through the Gamma-Poisson distribution, (2) the effects of confounding covariates can be effectively removed through a GLM framework, and (3) the DE genes show significant differences in means between case and control individuals after accounting for the individual-level covariates. Towards the first assumption, our posterior distribution effectively models counts through a Negative Binomial distribution, which has been justified for modeling scRNA-seq data [22] and has served as the foundation for many methods [13, 23, 26]. Towards the second assumption, many methods from the bulk RNA-seq data such as edgeR [9] and DESeq2 [8] have found tremendous success regressing out covariates through the GLM framework. Towards the last assumption, we note that there are existing methods such as IDEAS [27], scDD [28], and waddR [29] that test for differential *distributions* instead of differential mean expressions. However, we have found it more challenging to generalize these results when comparing different datasets of similar biological systems, as differentially distributed genes are not neatly characterized by over-/under-expression.

Design of simulation studies

The main message we wish to convey in this simulation section is two-fold: (1) testing for differential expression among individuals is fundamentally different from among cells, and (2) eSVD-DE's framework enables more accurate inference due to its usage of dimension reduction and posterior correction. Figure 2A conceptualizes the first point. Large within-individual variability hinders "pseudobulk" DE methods that sum the gene expressions across all the cells originating from the same individuals since these methods do not account for the variability of expression within an individual. On the other



Fig. 2 A Illustration of challenges for cohort-wide DE testing. **B** Setup for our simulation setup. **C** Isomap of the cells based on the true DE genes' expression before introducing the confounding variables. No individuals concentrate tightly in any region on the Isomap manifold, and there is a strong separation between the cases (shades of red) and controls (shades of blue). **D** Isomap of the observed data based on all the genes. Cells from the same individual concentrate in the embedding, suggesting that confounding covariates additionally drive the difference in expression profiles among individuals. **E** Downsampling experiment, demonstrating that by pooling information across genes, eSVD-DE outperforms gene-by-gene Negative Binomial regression for regressing out covariate effects. **F** Illustration of the importance of shrinkage, where the x-axis and y-axis represent each gene's test statistic with and without posterior correction, respectively. The genes are colored by their true log-fold change, of which the circled genes denote the top 50 genes with the highest true log-fold change. **G** ROC curve comparing four different methods, illustrating that eSVD-DE has more power than competing methods. The area under the curve (AUC) is shown for each method, where the percentage represents the area between the method's curve and the diagonal line as a fraction of total possible area. The bolded method denotes the method with highest AUC

hand, large between-individual variability hinders most existing methods designed for testing DE genes among cells from scRNA-seq data. This is because even if the difference in mean expression is insignificant on the cohort level, many cells from a small subset of individuals can yield highly significant differences in mean expression on the cell level. Since previous work has explained the importance of accounting for within-individual variability [27, 30, 31], we focus on a simulation that illustrates the importance of accounting for the between-individual variability here.

Our simulation study contains 700 genes and 20 individuals equally split among cases and controls, where each individual contributes 250 cells. The cells have gene expression profiles that are impacted by covariates correlated with the case-control status (such as age, sex, and tobacco use) and are sampled from a Gamma-Poisson distribution (Fig. 2B). Overall, the true generative model is more complex than the statistical model assumed by eSVD-DE to not give eSVD-DE an unfair advantage in our benchmarking experiment. Nonetheless, the data is generated such that 50 genes are drastically more differentially expressed among cases and controls on the cohort level than the remaining 650 genes, barring confounding effects. This true signal can be visualized via an Isomap [32, 33], where there is no apparent stratification within the case or control individuals (Fig. 2C). However, when the covariate effects of age, sex, and tobacco use are also included, there is obvious confounding of which genes are differentially expressed (Fig. 2D). Specifically, the cells from the same individual concentrate in different regions of the embedding. We choose to use the Isomap here, as opposed to the more commonly used UMAP [34], since the Isomap has well-studied statistical properties [35-37]. This quality is valuable when we use these visualizations to diagnose our hypothesis testing framework instead of an exploratory tool, as we will illustrate later.

Simulations verify that eSVD-DE effectively remove covariate effects in the presence of sparsity

We first show that by pooling information across genes via the dimension-reduction framework, eSVD-DE removes the confounding variables' effects more accurately compared to a Negative Binomial (NB) regression applied separately for each gene (Fig. 2E). To demonstrate this, we incrementally downsample the data to reduce the signal size of the DE genes. As the amount of downsampling increases, it will be more difficult for a method to remove the confounding effects properly. We measure the accuracy of how well the confounding effects were removed by computing the Pearson correlation between each cell's vector of true natural parameters across the 700 genes and its vector of estimated natural parameters after the confounding effects have been removed. We observe that at starting at 30% downsampling, eSVD-DE removes the confounding effects more accurately than the gene-by-gene NB regression (0.95-0.87). Additionally, although both methods drop in accuracy for downsampling levels of 60% or more, eSVD-DE is still more accurate than the gene-by-gene NB regression. This finding is statistically intuitive since when genes are sparsely observed, there is not enough information within a specific gene to accurately estimate the NB regression coefficients. However, by pooling information across genes, the coefficients for the confounding variables are better estimated.

Simulations verify that eSVD-DE's posterior enables gene-specific discoveries

We next illustrate that the posterior correction performed by eSVD-DE is important for DE testing on the cohort level (Fig. 2F). We plot the test statistic derived directly from the low-dimensional embedding against the test statistic derived from the posterior-corrected relative expression, where we mark the 50 genes the largest true log-fold change. In particular, the former represents the prototypical analysis of applying DE after denoising the data by a dimension-reduction method, and we observe that many null genes display large test statistics (x-axis), confirming the findings of previous work that observes an inflation of Type-1 error for such procedures [15, 30, 38]. However, by adjusting the relative expressions using the posterior mean, many null genes have drastically smaller test statistics (y-axis). This phenomenon occurs because the assumed linear model between the covariates and the gene's natural parameter does not sufficiently capture more complex non-parametric relationships often displayed among individuals, distorting the lower-dimensional embedding of the "residuals." This distortion has an adverse effect when genes vary substantially in sparsity—sparsely observed genes are denoised by projecting the cells onto the incorrect manifold.

Simulations verify that eSVD-DE's low-dimensional embedding improves power over other methods

Lastly, we illustrate the eSVD-DE has more power than other conventional methods to test for DE in cohort-wide scRNA-seq data through our simulation. Since different methods estimate sets of DE genes with dramatically different sizes for a particular FDR cutoff, we plot the entire curve of true positive rate (TPR) and false positive rate (FPR) over all possible cutoffs (Fig. 2G). Here, we compare against three other methods: DESeq2 [8] (i.e., the prototypical method representing "pseudobulk" analyses), MAST [10] (i.e., the commonly used method using mixed-effect models where there is a random effect among individuals) and SCTransform [39] (i.e., the prototypical method of performing DE ignoring the individual structure). We observe that eSVD-DE has the highest power compared to the three other methods. DESeq2 performs the best among the three competing methods since the averaging among cells of an individual dramatically reduces the estimation variability, while MAST performs the next best since it accounts for the individual structure.

We perform further power analyses across different number of genes and imbalances between number of cells across individuals (Additional file 1: Figs. C4 and C5). We also perform a separate simulation in the Appendix (Additional file 1: Figs. B2 and B3) to demonstrate that eSVD-DE does not inflate the Type-1 error among true null genes.

eSVD-DE enables diagnostics to assess the performance of removing covariate effects

Moving beyond simulations, we now investigate how well eSVD-DE removes confounding effects in scRNA-seq datasets and how the dimension-reduction framework enables conventional model diagnostics. To demonstrate this, we investigate a broad collection of scRNA-seq datasets from diverse tissues but focus here on the 8909 T-cells from a dataset of lung cells containing 10 healthy individuals and 24 individuals with IPF sequenced using the 10x Chromium single-cell platform, henceforth

Dataset	p	n	# Of case indiv.	# Of control indiv.	
Adams	6969	8909	24	10	
Habermann	6969	5286	6	4	
Smillie: TA 1 (I)	5713	6380	13	6	
Smillie: TA 1 (NI)	5713	14,215	13	6	
Velmeshev: L2/3	7094	12,984	15	16	

Table 1 Summary table of datase	ts
---------------------------------	----

called the "Adams" dataset [6]. We focus on 6969 genes for our analysis, including 5000 highly-variable genes, the genes reported to be DE by the authors as well as the 1003 human housekeeping genes [40]. We report the summary table of the Adams dataset, as well as the other datasets in this paper, in Table 1, where p denotes the number of genes and n denotes the number of cells. We include an extended summary table of the datasets in the Appendix When visualizing these cells via Isomap, we can see a clear separation of the case and control individuals, suggesting there are many genes with separable expression patterns (Fig. 3A). However, we also see that individual-level covariates like sex and smoking status also are locally concentrated in different regions of the Isomap. We do not wish to report genes as differentially expressed if the differences are induced only by sex differences or smoking.

After applying eSVD-DE, we see that the resulting Isomap of the fitted embedding no longer carries expression patterns correlated with the individuals, sex, or smoking status (Fig. 3B). If biologists prefer quantitative diagnostics, the practitioner can purposely omit a small percentage of values from the scRNA-seq count matrix before applying eSVD-DE and assess how correlated the predicted values are to these omitted values. This strategy was used successfully in our previous work [18]. For all these reasons, we advocate the matrix-factorization strategy in the eSVD-DE as opposed to deep-learning alternatives where there are fewer interpretable diagnostics available. The remaining diagnostics and the Isomaps for other scRNA-seq datasets are shown in the Appendix (Additional file 1: Fig. F6).

eSVD-DE appropriately adjusts for the sequencing depth using a dimension-reduction approach

Many authors have concluded that appropriately adjusting for the sequencing depth of each cell is one of the most influential aspects of an effective DE method [41, 42]. This adjustment is critical for scRNA-seq data since we do not wish to deem genes as DE simply because specific cells have a larger sequencing depth than others. Instead, genes should be deemed as DE if the case and controls have significantly different expression relative to the cells' sequencing depth. Accounting for this sequencing depth has been the source of numerous debates on how to normalize the cells' expression best [43]. The most commonly used approach is to log-normalize the gene expressions, but many existing work has cited that this normalization distorts qualitative aspects of the scRNA-seq dataset [12, 39].



Fig. 3 A Isomap of the leading principal components for the T-cells in the Adams dataset, shown by individuals, sex, and smoking status. The number in the insets denotes how correlated the individual-level covariate is correlated with the leading principal components. **B** Isomap of the cell embedding after applying eSVD-DE, demonstrating that confounding covariate and individual effects have been removed. The number in the insets denotes how correlated the individual-level covariate is correlated with the estimated eSVD-DE cell embedding. **C**, **D** Relationship between gene expression and cells' log sequencing depth before or after eSVD-DE (shown on a log scale), respectively, for groups of 6 genes partitioned by the gene's mean expression relationship after eSVD-DE. After eSVD-DE, each gene becomes more uncorrelated with the cell's log sequencing depth. **E** Relationship between each cell's log sequencing depth and overdispersion parameter, highlighting two genes displaying different shrinkage levels via the posterior mean

In order to quantify how well the sequencing depth effect was removed from the scRNA-seq data, we perform diagnostics analogous to those in [39]. Specifically, we group the genes into 6 bins based on their mean gene expression. We observe a significant relationship between each gene's expression and the cells' sequencing depth before normalizing the scRNA-seq data (Fig. 3C). This means without any accounting of the sequencing depth, genes could be deemed significantly differentially expressed solely due to cells having different sequencing depths. However, after fitting the eSVD-DE, this relationship is mainly removed across all bins, discounting the genes with the smallest mean expressions (Fig. 3D). Observe that eSVD-DE bears an advantage over other normalization methods such as SCTransform [39] and Scran [44] since eSVD-DE assesses the appropriate sequencing-depth normalization by

simultaneously accounting for the cells' gene expression and covariate information through its dimension-reduction framework, an important quality when dealing with cohort data. See Additional file 1: Fig. F7 for the application of this diagnostic on other datasets, and Additional file 1: Fig. F8 for how this diagnostic performs when only log-normalization is performed instead of the eSVD-DE in the Appendix.

Lastly, recall that our eSVD-DE framework adjusts the fitted gene expression values by the posterior Negative Binomial distribution. This adjustment relies on first quantifying the overdispersion of each gene, i.e., a higher overdispersion means the gene's expression patterns conform less to the estimated low-dimensional embedding. We hypothesize a negative correlation between the sequencing-depth normalization and the amount of overdispersion. This hypothesis is in line with previous works that investigate this relationship [39], citing that a smaller overdispersion implies that the fitted values are a good approximation of the gene expressions, which often occurs for densely observed genes. This can be visualized as a scatterplot (Fig. 3E). Additionally, this plot enables practitioners to survey the amount of shrinkage across all the genes broadly. We highlight RPS24 and SMAD3, 2 genes implicated in previous studies [45–47], as example genes that are highly and lowly shrunk via the posterior distribution since these genes were. All these aspects collectively support the claim that eSVD-DE is appropriately adjusting each gene's expression by the sequencing depth, which enables us to investigate the performance of the DE analyses downstream.

eSVD-DE recovers reproducible differences between case and control expression across multiple datasets

While it is difficult to assess the validity of DE genes in cohort-wide scRNA-seq data due to the lack of negative control genes, we hypothesize that a reliable proxy to assess the quality of the DE method is to collect two different cohort-wide scRNA-seq datasets of the same system and diseases/disorder and see if genes reported to be significant in one dataset display similar significances in the other. Towards this end, we investigate another dataset of 5286 T-cells of lung cells from 4 healthy individuals and 6 individuals with IPF sequenced using the 10x Chromium single-cell platform of the same 6969 genes here, henceforth called the "Habermann" dataset (Fig. 4A) [19]. Since we are primarily interested in comparing eSVD-DE to other DE methods using this pair of datasets, we do not deploy a multiple-testing procedure to select DE genes but investigate the sets genes with the largest test statistic magnitudes of the same cardinality as those reported by the authors for meaningful comparisons. For starters, the volcano plot on the Adams dataset shows that 20 of the 84 genes with the largest test statistics derived from our eSVD-DE procedure intersect with the 84 DE genes reported by the authors, resulting in Fisher's exact test *p*-value of 2.2×10^{-21} (Fig. 4B). Similarly, the volcano plot on the Habermann dataset shows that 30 of the 157 genes with the largest test statistics derived from our eSVD-DE procedure intersect with the 157 DE genes reported by the authors, resulting in Fisher's exact test *p*-value of 3.1×10^{-20} (Fig. 4C). Additionally, since housekeeping genes constitute genes primarily responsible for basic cellular functions and are stably expressed regardless of cellular condition, we hypothesize that these genes should not carry substantial differential expression patterns [40]. This phenomenon is demonstrated through the volcano plots (Fig. 4B, C).



Fig. 4 A Number of case and control individuals and the total number of cells across cell types and datasets. The x-axis denotes the total number of cells on a log scale, while the partitioning of case and control cells denotes the number of cells from each individual relative to the total number of cells. **B**, **C**, **G**, **H** Volcano plot, where the set of genes with the largest test statistics (having the same size as the original author's set) is depicted in orange. The genes reported by the respective authors or housekeeping genes are in purple and green, respectively. The Fisher exact test's *p*-value between the enrichment of eSVD-DE's DE genes and the author's reported DE genes is also reported. **D** Upset plot showing the intersection between pairs of DE genes, either the reported DE genes in the Habermann and Adams dataset and the estimated DE genes using eSVD-DE or DESeq2. **E**, **I** Hexplots showing the correlation between genes' eSVD-DE test statistics across datasets, either the originally reported DE genes in either dataset or the housekeeping genes. **F**, **J** Similar to **E**, **I** but showing the gene's DESeq2 test statistics

When comparing different DE methods, we ask if the leading genes derived from the Adams dataset using eSVD-DE are either (1) genes reported by authors of the Habermann dataset or (2) themselves the leading genes derived from the Habermann dataset using eSVD-DE, or vice-versa. If the size of this intersection is larger than those derived by other DE methods in either scenario, we would have partial evidence that eSVD-DE is more consistently recovering reproducible signals across both datasets. We compared

against DESeq2 whereby cells of each individual are aggregated into "pseudobulk" expressions prior to performing DESeq2 [8] for this investigation. We highlight this pseudobulk procedure specifically since this procedure was reported to be more reliable for scRNA-seq with biologically-replicate samples, but its reliability for cohort data was not investigated [4]. For this comparison, we select the 226 genes with the largest test statistics in magnitude from either the Adams or Habermann dataset using either eSVD-DE or DESeq2 since the 84 reported DE genes for the Adams dataset and 157 reported DE genes for the Habermann dataset yield 226 unique genes. The upset plot shows that eSVD-DE yields larger intersections between the two datasets than DESeq2 (Fig. 4D).

Additionally, we hypothesize that if a DE method recovers reproducible signals, the DE genes test statistics should be positively correlated between the two datasets. Indeed, if a reported DE gene has a positive log-fold change in one dataset, the other dataset should also show evidence of a positive change. On the other hand, we hypothesize that genes unrelated to the disease/disorder should be uncorrelated test statistics between the two datasets. This hypothesis would be biologically justifiable, as the log-fold change of a gene unrelated to disease/disorder would be determined by random chance. We observe these relationships for eSVD-DE's test statistics (Fig. 4E). In contrast, the correlation for DESeq2's test statistic among the reported DE genes is near-zero (Fig. 4F). We suspect all the above phenomenons are likely driven by pseudobulk methods' lack of accounting for the within-individual variability.

To further support our empirical claims regarding eSVD-DE, we also study cells from 18 individuals with ulcerative colitis (UC) and 12 healthy individuals across multiple cell types in the colon and 5713 genes, henceforth called the "Smillie" dataset [20]. Each individual with UC contributed cells from inflamed and non-inflamed colon biopsies, and each healthy individual contributed two biologically replicated samples from biopsies in analogous colon regions. In our analysis, we treat each tissue sample as a different "individual" and perform two DE analyses—one of non-inflamed samples from individuals with UC against healthy samples (i.e., the "non-inflamed" analysis) and another of inflamed samples from individuals with UC against healthy samples (i.e., the "inflamed" analysis). Importantly, we split the healthy tissue samples so each healthy sample is only involved in one of the two DE analyses. While we apply this pair of analyses on multiple cell types in this biological system, we focus on the analysis of transit amplifying 1 (TA 1) cells here (Fig. 4A). Similar to before, we see that the volcano plots for both the non-inflamed and inflamed DE analysis yield highly-enriched intersections between the genes with the largest test statistics and the DE genes reported by the authors, as well as near-zero enrichment of the housekeeping genes (Fig. 4G,H).

The authors of the Smillie dataset observed a high correlation among the test statistics of reported DE genes between the non-inflamed and inflamed DE analyses when aggregating across all the cell types, suggesting that the transcriptomic signature of UC precedes inflammation [20]. We hypothesize that a higher-powered DE analysis of cohort-wide scRNA data should reveal a strong correlation even when focusing on only one cell type. Indeed, we see this positive correlation among eSVD-DE's test statistics ($\rho = 0.84$, Fig. 4I). Additionally, while we see a positive correlation among the housekeeping genes ($\rho = 0.49$), many of such genes have a near-zero test statistic in the non-inflamed analysis. This observation suggests that the differences in housekeeping genes' expression are induced more by the inflammation of the tissue rather than a biological mechanism disrupted by ulcerative colitis. In contrast, when we apply a similar pair of analyses using DESeq2 on pseudobulk data, the correlation among the reported DE genes is substantially lower (Fig. 4J). See Additional file 2 for the resulting statistics when analyzing the Adams and Habermann dataset, and Additional file 3 for the Smillie datasets.

eSVD-DE detects DE genes highly enriched with previously-annotated genes

Having demonstrated all the advantages of eSVD-DE, we hypothesize that our method leads to novel cell-type specific DE discovery. Towards this end, we analyze brain



Fig. 5 A Number of case and control individuals and the total number of cells across cell types for the Velmeshev dataset, analogous to Fig. 4A. **B** Volcano plot showing eSVD-DE's results for layer 2/3 cells, where the set of genes exceeding the empirical FDR cutoff are depicted in orange. The SFARI, bulk DE, and housekeeping genes are in purple, blue, and green, respectively. The Fisher exact *p*-value between the selected genes via eSVD-DE and the bulk DE genes is noted. **C** Comparison of the eSVD-DE *p*-values against the DESeq2 *p*-values, and the genes exceeding an FDR cutoff for DESeq2 are shown in yellow. The correlation between the two sets of negative log₁₀ *p*-values is noted. **D** Increase in enrichment for relevant GO terms among the 331 DE genes found by eSVD-DE, the 144 DE genes found by DESeq2, and the 109 genes initially reported by the authors. **E** Upset plot comparing the SFARI or bulk DE genes with the top 100 DE genes estimated by eSVD-DE, DESeq2, MAST, or SCTransform. **F** Downsampling experiment, illustrating the stability of the 50 genes with the largest test statistics in magnitude among the SFARI or the bulk DE genes as the dataset is artificially downsampled. The values on the *y*-axis is the ratio of the test statistic between these genes and the housekeeping genes

single-cells from controls and case individuals who had been diagnosed with autism spectrum disorder (ASD), partitioned among many cell-types [21], henceforth called the "Velmeshev" dataset (Fig. 5A). We focus on this particular system since the Simons Foundation Autism Research Initiative (SFARI) routinely curates a list of autism risk genes based on genetic studies from many publications [48], which can provide relevant genes to compare our eSVD-DE results against. Additionally, a recent bulk-DE analysis of ASD provides a high-quality and independent list of genes to compare against, which we will call the "bulk DE genes" henceforth [49]. A priori, we would not expect any method deployed on the Velmeshev dataset to fully match these lists because our analysis is cell-type specific and the SFARI list are genetic risk genes, which do not always lead to differential expression.

When analyzing the 12,984 cells in layer 2/3, we observe that eSVD-DE estimates 331 DE genes (among the 7055 genes used in the analysis) using an empirical FDR cutoff of 0.05, where we calibrate the *p*-value according to an empirical null distribution to account for potential misspecification [50] (Fig. 5B). Here, 42 and 95 of these genes are among the SFARI and the bulk DE genes, respectively (among the total 800 and 1556 genes, respectively). Comparing the bulk DE genes to our eSVD-DE genes (which are both derived from transcriptomics data), we obtain a significant Fisher p-value of 0.002. Additionally, when we compare DESeq2 for the same cells to find DE genes (144), only 19 and 30 genes overlapped with the SFARI and the bulk DE genes, respectively (Fig. 5C). This finding demonstrates that eSVD-DE can find more DE genes, which are more relevant, compared to DESeq2. We report analogous plots in the Appendix when comparing against other methods, such as MAST and SCTransform (Additional file 1: Fig. F9). We note that we do not compare against the DE genes found by the authors of the data themselves here, as they had used MAST to find their DE genes. As our simulations beforehand demonstrated, MAST is not necessarily a reliable "gold standard" to compare against. However, the 331 DE genes found by eSVD-DE are also much more enriched in GO terms that are plausibly related to ASD when compared to the DESeq2 genes or the original genes reported for layer 2/3 (Fig. 5D). All these observations suggest that eSVD-DE is well-equipped to uncover meaningful results.

Next, we investigated qualitative differences among the methods compared to the SFARI and bulk DE genes. We enable fair comparisons among the different methods by finding the top 100 genes with the smallest *p*-values for each method. We then counted how many of these 100 genes intersected with the SFARI or bulk DE genes (Fig. 5E). eSVD-DE has the highest overlap with both sets (31 with the bulk DE, 13 with the SFARI genes), followed closely by DESeq2 and then finally by MAST and SCTransform. In fact, eSVD-DE and DESeq2 have an overlap of 27 genes. The observations from Fig. 5D, E combined suggest that eSVD-DE inherits many advantages of pseudobulk methods while improving these methods by accounting for within-individual variability. In the second investigation, we asked how stable each method was as the dataset was gradually downsampled. Indeed, as the scRNA-seq dataset becomes sparser, the relative difference between "strongly-expressed" DE genes and null genes becomes fainter, and we would expect any method to degrade in performance. To investigate this, we take the 50 genes with the largest test statistics in magnitude among the SFARI and bulk DE genes for each method and ask how

their mean test statistic (relative to the housekeeping gene's test statistic) diminishes as the scRNA-seq data is downsampled. We use the housekeeping genes as a proxy for the null genes, which is a reasonable choice given Fig. 5B. As the downsampling percentage increases (i.e., the signal becomes fainter), MAST and SCTransform lose their ability to distinguish between formally highly-differentiable genes and housekeeping genes. This finding makes sense, as MAST and SCTransform estimate the DE genes based on a gene-by-gene regression, meaning no information is shared between genes. On the other hand, DESeq2 and eSVD-DE are surprisingly stable—they can demonstrate a clear separation between the 50 DE genes and the housekeeping genes even when the data is downsampled at 40%. This finding is also sensible, as DESeq2 aggregates cells among individuals, and eSVD-DE pools information between genes; both strategies yield robustness against sparsity. However, eSVD-DE is still preferred over DESeq2 here, as the separation between the DE genes and housekeeping genes is much higher for eSVD-DE.

We include additional diagnostic plots via Isomaps, volcano plots of the DE genes, and the GO analysis results in Appendix (Additional file 1: Fig. F7, and F10 through F12). See Additional file 4 for the resulting statistics when analyzing the Velmeshev datasets.

Conclusions

We demonstrate the nuances of performing cohort-wide differential testing for singlecell RNA-seq data. The difficulty primarily stems frotm individual-level confounding covariates, which can be difficult to remove using typical DE strategies of regressing out their effects gene-by-gene. Instead, eSVD-DE pools the information across genes to remove the confounding covariates, yielding empirical performance that is more promising than current pseudobulk DE methods or DE methods currently used for scRNA-seq data when no prevalent individual-level covariates are present. We achieve this through a matrix factorization strategy, estimating the coefficients associated with the covariates and each cell's and gene's latent vectors. We then shrink the estimated denoised gene expressions via the posterior mean according to the Gamma-Poisson distribution. This strategy helps dampen potential over-smoothing effects induced by the matrix factorization. We then deploy a test statistic designed to test for DE genes on the individual level instead of the cellular level. Note that our procedure can be used concurrently with IDEAS [27], which aims to find differential distributions between cases and controls, which could be more challenging to interpret than differential means. However, we do not pursue this direction in this paper. We also note that gene selection for cohort-level scRNA-seq datasets is an important task that we are interested in exploring in future work, since the inclusion of certain genes could impact the *p*-value of other genes due to the nature of pooling information across cells and genes via a low-dimensional embedding. Ideas from graph-representation work such as [51] and [52] could be highly relevant in this direction. We hope that eSVD-DE would be beneficial for inspiring cohort-wide DE tests for future single-cell assays beyond scRNA-seq and single-cell eQTL analyses where abundant individual-level covariate effects must be adequately removed. Additionally, we are curious about broader settings where instead of having case and control individuals within a cohort, we are interested in testing if continuous covariates such as age have a substantial transcriptomic impact on specific cell types, as discussed in [53].

Statistical model and method

Let $A \in \{0, 1, ..., \}^{p \times n}$ denote the observed count matrix with *n* cells and *p* genes, and $C \in \mathbb{R}^{n \times r}$ denote the observed *r* covariates for the *n* cells. Importantly, certain columns of *C* would be the following:

- **Intercept**: Let $C_{.,1} = 1$. We'll call this column $C_{.,(int)}$.
- Log sequencing depth: Let $s_i = \sum_{j=1}^p A_{ji}$. Then, let $C_{i,2} = \log(s_i)$ for all $i \in \{1, ..., n\}$. We'll call this column $C_{\cdot,(\text{lib})}$.
- **Case–control status**: Let $C_{i,3} \in \{0, 1\}$ depending on whether or not cell $i \in \{1, ..., n\}$ is associated with a case or control individual. We'll call this column $C_{i,(cc)}$.
- Others: The remaining columns of *C* could contain numerical covariates associated with the cells: which region of the body the cell was sampled from, the age of the corresponding individual, etc. We assume the categorical covariates have already been transformed via one-hot encodings (i.e., categorical variables with *k* levels are transformed using k 1 indicator variables). In contrast, assume the numerical covariates are standardized to have a standard deviation of 1. We have found it beneficial to not center the numerical covariates around 0. This way, the resulting estimated coefficients in *Z* are easier to interpret.

Optionally, practitioners may include one-hot encoding vectors for which cells originate from which individuals. In our paper, we have found this to be optional and sometimes detrimental to the fit. This is because the inclusion of such one-hot encoding vectors result in a collinear matrix *C*. This preparation of the covariate matrix *C* is primarily handled by the function eSVD2::.reparameterization_esvd_covariates in our codebase.

The statistical foundation of eSVD-DE is the following hierarchical model where each entry of A_{ji} is modeled as,

$$(A_{ji}|\lambda_{ji}) \sim \text{Poisson}(\ell_{ji} \cdot \lambda_{ji}), \text{ and } \lambda_{ji} \sim \text{Gamma}(\alpha = \mu_{ji}/\gamma_j; \beta = 1/\gamma_j),$$
 (1)

(where α and β denote the shape and rate parameters of a Gamma distribution respectively) for gene $j \in \{1, ..., p\}$ and cell $i \in \{1, ..., n\}$, where the low-dimensional mean matrix is $\mu \in \mathbb{R}^{p \times n}$ where

$$\mu_{ji} = \exp\left((Y_{j,\cdot})^{\top}(X_{i,\cdot}) + Z_{j,(cc)} \cdot C_{i,(cc)}\right),\tag{2}$$

for $X \in \mathbb{R}^{n \times k}$ and $Y \in \mathbb{R}^{p \times k}$, and ℓ_{ji} denotes the covariate-adjusted sequencing depth,

$$\ell_{ji} = \exp\left(\left(Z_{j,-(\mathrm{cc})}\right)^{\top}C_{i,-(\mathrm{cc})}\right),\tag{3}$$

and $\gamma_j > 0$ denotes the overdispersion parameter for gene *j* which controls the variability of gene *j*. Here, the subset notation "–(cc)" means that we exclude the "(cc)" column from the corresponding matrix. In the terminology of previous work [22, 23], we

interpret μ_{ji} as the "predictable" gene expression (i.e., expression of gene *j* in cell *i* that can be predicted from other cells and genes, which we interpret as a low-dimensional manifold), and λ_{ji} as the biological relative expression of gene *j* in cell *i*. We observe the count matrix *A* as well as the covariates *C*, and we need to estimate the cells' latent embedding *X*, the genes' latent embedding *Y*, the covariates' coefficients *Z* as well as the overdispersion vector γ .

Regarding the Gamma distribution

Our parameterization of the Gamma distribution is inspired by the constant Fano factor model used in SAVER [23]. Specifically, for the Gamma distribution in (1),

$$\mathbb{E}[\lambda_{ji}] = \mu_{ji}$$
, and $\mathbb{V}[\lambda_{ji}] = \gamma_j \cdot \mu_{ji}$,

meaning the variance scales proportionally with the mean. Additionally, it can be derived that the marginal distribution of A_{ji} is a negative-binomial,

$$A_{ji} \sim \mathrm{NB}(r = \mu_{ji}/\gamma_j; \ p = \ell_{ji}/(\ell_{ji} + 1/\gamma_j)),$$

meaning the moments of A_{ji} marginally are

$$\mathbb{E}[A_{ji}] = \ell_{ji} \cdot \mu_{ji}, \text{ and } \mathbb{V}[A_{ji}] = (\ell_{ji}\gamma_j + 1) \cdot \ell_{ji} \cdot \mu_{ji}$$

These equations demonstrate that γ_j measures overdispersion—the larger γ_j is, the larger the variance of A_{ii} is.

Additionally, we can derive that the posterior distribution of $\lambda_{ji}|A_{ji}$ is

$$\lambda_{ji}|A_{ji} \sim \text{Gamma}(\alpha = A_{ji} + \mu_{ji}/\gamma_j; \ \beta = \ell_{ji} + 1/\gamma_j), \tag{4}$$

which will be useful when deriving our method below.

Rationale and justification of statistical model

Our hierarchical model (1) draws inspiration from two literatures. The first literature is the single-cell modeling literature. Various work have shown that hierarchical model (1) that model observed counts as a Poisson distribution originating from a Gamma prior is a statistically sound and empirically justified model of scRNA-seq data [22, 23]. These models are appealing as they explicitly separate the technical noise (i.e., noise incurred by the nature of sequencing, modeled via the Poisson distribution) from the biological noise (i.e., noise naturally occurring among cells, modeled via the Gamma distribution). Other work has provided additional biological explanations on why biological noise is commonplace in organisms [54].

The second literature is the differential-expression literature, which canonically focuses on testing for the differential expression of a gene, one gene at a time. For gene j, methods like DESeq2 [8], MAST [10], and SCTransform [39] regress the observed expression A_{j1}, \ldots, A_{jn} onto observed sequencing depth of each cell ℓ_1, \ldots, ℓ_n . These models account for the possibility that the sequencing depth detrimentally confounds the differential expression of gene j in varying degrees across all genes. Hence, our modeling of the sequencing depth in (3) is qualitatively similar.

High-level description of the hypothesis test

We provide a high-level description of the hypothesis that eSVD-DE tests, with respect to the hierarchical model prescribed in the aforementioned statistical model. Let \mathcal{A} denote the set of case individuals and \mathcal{B} denote the set of control individuals. Consider a gene $j \in \{1, ..., p\}$, assume that all there are parameters $\mu_j^{(\text{case})}$ and $v_j^{(\text{case})}$ that depict the unknown gene expression mean and variance among the case individuals, i.e.,

$$\mu_j^{(s)} \sim \text{Gaussian}(\mu_j^{(\text{case})}, \nu_j^{(\text{case})}), \text{ for each } s \in \mathcal{A},$$

and likewise for the control individuals for unknown parameters $\mu^{(\text{control})}$ and $\sigma^{(\text{control})}$,

$$\mu_j^{(s)} \sim \text{Gaussian}(\mu_j^{(\text{control})}, \nu_j^{(\text{control})}), \text{ for each } s \in \mathcal{B}.$$

Each individual then contributes multiple cells. Let $\mathcal{I}(s) \subset \{1, ..., n\}$ denote the set of cells that an individual contributes. We model cell *i*'s expression (from individual *s*, where *s* refers to "subject") for gene *j* as

$$\lambda_{ji} \sim F_j^{(s)}, \quad \text{for } i \in \mathcal{I}(s),$$
(5)

where $F_j^{(s)}$ is a distribution with mean $\mu_j^{(s)}$. From here, the observed count in cell *i* for gene *j*, defined as A_{ji} is related to λ_{ji} through (1).

With this alternative perspective of our statistical model, eSVD-DE is testing between the null and alternative hypotheses,

$$H_{0,j}: \mu_j^{\text{(case)}} = \mu_j^{\text{(control)}}, \quad \text{and} \quad H_{1,j}: \mu_j^{\text{(case)}} \neq \mu_j^{\text{(control)}}, \tag{6}$$

which are comparing the pouplation case-individual mean expression to the pouplation control-individual mean expression of gene *j*.

A few notes are in order:

 \sim

- Rationale of against pseudobulk approaches: The hypothesis test in (6) might suggest a pseudobulk approach (such as used by DESeq2 in our comparisons). However, as we have mentioned in the main text, such pseudobulk methods do not capture variability within an individual (i.e., the variance of the distributions $F_j^{(s)}$). This can affect validity of such pseudobulk strategies, as demonstrated in our null simulations shown later.
- **Rationale of our hypothesis test**: We note that our hypothesis test in (6) relates the population case individual's against the population control individual's mean expression. This is different from the following null hypotheses that are not as apt for testing for differential mean expression for cohort-wide scRNA-seq data:

$$H'_{0,j}: \operatorname{Mean}_{i \in \mathcal{I}(s), s \in \mathcal{A}}(\{\lambda_{ij}\}) = \operatorname{Mean}_{i \in \mathcal{I}(s), s \in \mathcal{B}}(\{\lambda_{ij}\}),$$

$$\tag{7}$$

or

$$H_{0,i}'': \operatorname{Mean}_{i \in \mathcal{I}(s), s \in \mathcal{A}}(\{\mu_{ij}\}) = \operatorname{Mean}_{i \in \mathcal{I}(s), s \in \mathcal{B}}(\{\mu_{ij}\}).$$

$$(8)$$

We cannot test the null hypothesis in (7) directly since we observe only one observation A_{ji} for each λ_{ji} . By pooling all the information across cells and genes, we can instead estimate μ_{ji} , which is the mean of λ_{ji} based on a low-rank matrix factorization. In contrast, (8) ignores which cells originated from which individual, which makes it undesirable for differential testing for cohort-wide scRNA-seq data. For example, if all the cells from a few case individuals have vastly higher gene expression than all other case individuals' cells, the null hypothesis in (8) might be rejected, but eSVD-DE's null hypothesis in (6) would not be rejected. In such a scenario, we would not want to reject the null hypothesis since the behavior of a small number of case individuals would not be representative of the human population of case individuals. In our comparisons, methods like SCTransform test the null hypothesis in (8). Methods like SCTransform face a different computational obstacle. Since these methods regress out the individuals' substantial covariate effects one gene at a time, this regression might be inaccurately estimated due to the sparsity of scRNA-seq data.

- **Distinction between cell-mean and individual-mean**: Observe that even though our hypothesis testing framework in (5) treats all the λ_{ji} 's for the same individual (i.e., $i \in \mathcal{I}(s)$) as i.i.d., eSVD-DE nonetheless models each cell's mean as μ_{ji} in (2). This is to facilitate the matrix factorization framework in order to pool information across cells and genes. As we will discuss later in the model, we estimate $\mu_{ji}^{(s)}$ by averaging the estimated μ_{ji} 's (after taking its posterior distribution to account for model misspecification, see (14)). We believe it is important to rely on the posterior distribution especially for cohort-wide differential expression testing since single-cell sequencing data is quite sparse, so it's important to have a data-driven procedure to adjust our estimated values of μ_{ji} that balances how sparsely sequenced (or said differently, how much "information") gene *j* is and how large the covariate-adjusted sequencing depth of cell *i* is.
- **Covariate effects, sparsity, and overdispersion**: Within this perspective of the model, all the confounding covariate effects is modeled through the covariate-adjusted sequencing depth ℓ_{ji} in (3), which eSVD-DE tries to model appropriately by pooling information across cells and genes. Additionally, we lose power to test for (6) as sparsity increases, which is parameterized by the overdispersion parameter γ_{ji} .

Implementation and details of eSVD-DE

Initialization

The initialization broadly falls into three steps: (1) initializing the estimate of the coefficients $\hat{Z} \in \mathbb{R}^{p \times r}$, (2) initializing the estimate of the embeddings for the cells $\hat{X} \in \mathbb{R}^{n \times k}$ and for the genes $\hat{Y} \in \mathbb{R}^{p \times k}$, and 3) reparameterizing our estimates \hat{X} , \hat{Y} and \hat{Z} . We describe each step below.

First, to initialize our estimate \hat{Z} , we fit a Poisson regression with a ridge penalty separately for each cell *i*. This regresses the covariates *C* onto *A*_{.,*i*}, the observed gene counts for cell *i*. Specifically, for a small penalty $\tau > 0$, we initialize the *j*th row of \hat{Z} as

$$\hat{Z}_{j,\cdot} = \operatorname*{argmin}_{z \in \mathbb{R}^r} - \left[\sum_{i=1}^n A_{ji} \cdot (z^\top C_{i,\cdot}) - \exp(z^\top C_{i,\cdot}) \right] + \tau \cdot \|z_{-1}\|_2^2, \quad \text{for } j \in \{1, \dots, p\},$$

where $z_{-1} = (z_2, ..., z_r) \in \mathbb{R}^{r-1}$. The ridge penalty omits z_1 because, by definition, $C_{\cdot,1}$ is the all-one vector that represents the intercept. Typically, we set $\tau = 0.01$, a small non-negative value to mitigate multicollinearity issues, similar to other works such as [12].

Second, we initialize our estimates \hat{X} and \hat{Y} . Consider the matrix,

$$R = \log(A+1) - \hat{Z}^{\top}C \in \mathbb{R}^{p \times n}$$

where $log(\cdot)$ of a matrix denotes taking the natural logarithm of each matrix entry. Consider the SVD of *R*,

$$R = UDV^{\top},$$

where $U \in \mathbb{R}^{p \times k}$ and $V \in \mathbb{R}^{n \times k}$ are column-wise orthonormal matrices, and D is a diagonal non-negative matrix with decreasing values along the diagonal. We then initialize our estimates of \hat{X} and \hat{Y} as

$$\hat{X} = V\sqrt{D} \in \mathbb{R}^{n \times k}$$
, and $\hat{Y} = U\sqrt{D} \in \mathbb{R}^{p \times k}$,

where $\sqrt{\cdot}$ of a diagonal non-negative matrix denotes taking the square root of all its diagonal entries. This initialization of \hat{X} and \hat{Y} is motivated by the observation that if we model A_{ji} as a Poisson random variable, then $\log(A_{ji} + 1)$ is a crude approximation of the natural parameter μ_{ji} , and $(\hat{Z}_{j,\cdot})^{\top}C_{i,\cdot}$ would be an initial estimate of the covariate effect on μ_{ji} .

This step is collectively handled by eSVD2::initialize_esvd in our codebase. It relies on the glmnet::glmnet function to fit the Poisson ridge regression.

Optimization of the embeddings

The iterative optimization of the embedding broadly falls into two steps to estimate the latent factors of the mean matrix μ defined in (2): (1) updating the estimates of \hat{X} , \hat{Y} , and \hat{Z} , holding the initialized values of $\hat{Z}_{.,(cc)}$ fixed, and then (2) updating the estimates of \hat{X} , \hat{Y} , and \hat{Z} , allowing the values in $\hat{Z}_{.,(cc)}$ to be updated. This step takes inspiration from methods such as GLM-PCA [12], ZINB-WaVE [13], and our previous work [18]. We first describe the details of the procedure, followed by its justification.

For a tuning parameter $\tau > 0$, we seek to optimize the following objective function

$$\{\hat{X}, \hat{Y}, \hat{Z}\} = \underset{X, Y, Z}{\operatorname{argmin}} \left[-\frac{1}{np} \left[\sum_{ij} \log P(A_{ji} \mid X_{i, \cdot}, Y_{j, \cdot}, C_{i, \cdot}, Z_{j, \cdot}) \right] + \tau \cdot \left[\|X\|_F^2 + \|Y\|_F^2 + \|Z\|_F^2 \right].$$
(9)

Let us define the natural parameter,

$$\theta_{ji} = (Y_{j,\cdot})^{\top} X_{i,\cdot} + (Z_{j,\cdot})^{\top} C_{i,\cdot}, \text{ for all } i \in \{1, \ldots, n\}, j \in \{1, \ldots, p\}.$$

Then, here, $\log P(A_{ji} | X_{i,\cdot}, Y_{j,\cdot}, C_{i,\cdot}, Z_{j,\cdot})$ denotes the log-likelihood for a Poisson random variable,

$$\log P(A_{ji} \mid X_{i,\cdot}, Y_{j,\cdot}, C_{i,\cdot}, Z_{j,\cdot}) = A_{ji} \cdot \theta_{ji} - \exp(\theta_{ji}).$$

The objective function in (9) is a non-convex, as documented in [18]. This poses challenging optimization considerations. Hence, we first describe an alternating minimization strategy that improves upon previous work computationally. Then, we describe the aforementioned two-step approach.

Our alternating minimization consists of the following steps, starting on an initial estimation of $\hat{X}^{(0)} = \hat{X}$, $\hat{Y}^{(0)} = \hat{Y}$, and $\hat{Z}^{(0)} = \hat{Z}$. This strategy is motivated by the observation that holding *Y* and *Z* fixed, the optimization in (9) over *X* is convex, and vice-versa. Let *t* denote the iteration counter.

• Optimize the cell embedding,

$$\hat{X}^{(t+1)} = \underset{X}{\operatorname{argmin}} \left[-\frac{1}{np} \left[\sum_{ij} \log P(A_{ji} \mid X_{i,\cdot}, \hat{Y}^{(t)}_{j,\cdot}, C_{i,\cdot}, \hat{Z}^{(t)}_{j,\cdot}) \right] + \tau \cdot \left[\|X\|_F^2 \right].$$
(10)

· Optimize the gene embedding and coefficients

$$\left\{\hat{Y}^{(t+1)}, \hat{Z}^{(t+1)}\right\} = \underset{Y,Z}{\operatorname{argmin}} \left[-\frac{1}{np} \left[\sum_{ij} \log P(A_{ji} \mid \hat{X}^{(t+1)}_{i,\cdot}, Y_{j,\cdot}, C_{i,\cdot}, Z_{j,\cdot})\right] + \tau \cdot \left[\|Y\|_F^2 + \|Z\|_F^2\right].$$
(11)

We repeat until convergence and then perform a reparameterization (described in the next section).

The optimizations (10) and (11) are performed using a Newton optimization, which is a second-order method. This specific optimization framework is ideal for solving (10) and (11) for a few reasons:

1. **Parallelization into many low-dimension optimizations**: Both optimizations (10) and (11) decompose into *n* and *p* smaller optimization problems respectively. For example, the *i*th row of $\hat{X}^{(t+1)}$ is equivalently solved by the optimization

$$\underset{x \in \mathbb{R}^{K}}{\operatorname{argmin}} \left[-\frac{1}{np} \sum_{j=1}^{p} \log P(A_{ji} \mid x, \hat{Y}_{j, \cdot}^{(t)}, C_{i, \cdot}, \hat{Z}_{j, \cdot}^{(t)}) \right] + \tau \cdot \|x\|_{2}^{2}.$$
(12)

Hence, solving (10) and (11) amounts to a solving *n* different *K*-dimensional or *p* different (K + r)-dimensional optimizations respectively. Since we are in a setting where $\max(K, r) \ll \min(n, p)$, it is beneficial to use a second-order method since the Hessian information can yield faster convergence rates in terms of the number of iterations needed to solve optimizations like (12), and there is not a large computational overhead to compute the Hessians to solve (12) (especially when *P* is the Poisson distribution, where the Hessians are straight-forward to derive and compute).

2. No randomness in the optimization: Since (10) and (11) decompose into n and p smaller optimization problems respectively, these optimizations can be embarrassingly parallelized. Hence, there is no need to consider stochastic optimization schemes. This is appealing as this means different practitioners using our method would necessarily obtain the same resulting fit.

3. Generalization to other exponential families: While we focus on specifically solving (10) and (11) for the Poisson distribution, using a Newton optimization is ideal for other exponential-family distributions as well. This is because while other exponential-family distributions like the exponential and Negative Binomial have constraints on the natural parameters θ_{ij} , the gradient and Hessians of these log-likelihoods naturally prevent the optimization iterates from violating these constraints. Hence, our codebase can handle modeling situations beyond this paper's scope.

Our optimization procedure is then the following. First, starting from the current estimates of \hat{X} , \hat{Y} , and \hat{Z} in the previous step ("Initialization"). we updating the estimates of \hat{X} , \hat{Y} , and \hat{Z} via alternating minimization holding the initialized values of $\hat{Z}_{.,(cc)}$ fixed (what we'll call the "Phase one optimization"). After convergence, we then perform a second round of alternating minimization to update the estimates of \hat{X} , \hat{Y} , and \hat{Z} , allowing the values in $\hat{Z}_{.,(cc)}$ to be updated (what we'll call the "Phase two optimization"). (Here, we omit the superscript "(*t*)" for notational simplicity.) We do these two phases of optimization since empirically, we have observed more stable behavior and a better final objective value doing two rounds of optimization (compared to optimizing all of *X*, *Y*, and *Z*, including $Z_{.,(cc)}$ from the start). This becomes imperative since $Z_{.,(cc)}$ will play a more pronounced role in the remainder of eSVD-DE compared to all the other covariates in *Z*, because is it part of the "signal" that we wish to estimate and is not a "confounder." This is inspired by theoretical results regarding warm-starting the non-convex optimization [55, 56]. This step is collectively handled by eSVD2::opt_esvd in our codebase.

Reparameterizing the matrix factorization

After completing either of the two phases mentioned above, we apply the following reparameterization procedure. The goal of this reparameterization is to ensure identifiability since, a priori, there could be multiple estimates $\{\hat{X}, \hat{Y}, \hat{Z}\}$ that have the same predictive power but offer different interpretations of the data. For instance, without such a reparameterization, there could be many columns in \hat{X} that are correlated with other columns in \hat{X} or columns in the covariate matrix *C*, which would obfuscate interpreting different axes of variation. Hence, our reparameterization procedure ensures orthogonality among our estimated matrix factorization to resolve this identifiability concern.

When describing our reparameterization procedure, for notational simplicity, we let $\{\hat{X}, \hat{Y}, \hat{Z}\}$ denote $\{\hat{X}^{(T)}, \hat{Y}^{(T)}, \hat{Z}^{(T)}\}$, which is the final estimate after *T* iterations for either the Phase one or two optimizations. We also let the " $a \leftarrow b$ " notation denote setting the variable *a* to be the value in variable *b*.

The reparameterization procedure operates in two steps. The first step ensures that \hat{X} is orthogonal to *C* (which would result in adjusting \hat{X} and \hat{Z}). This ensures that \hat{X} captures variability that is not explained by *C*. The second step ensures that both \hat{X} and \hat{Y} have orthogonal columns (which would result in adjusting \hat{X} and \hat{Y}). This ensures that each respective latent dimension in both \hat{X} and \hat{Y} are identifiable.

• **Step 1**: We first perform a regression of each latent dimension of \hat{X} onto *C*. That is, for each latent dimension $d \in \{1, ..., k\}$,

$$\hat{X}_{\cdot,d} = C^{\top}\beta + \epsilon,$$

where $\beta \in \mathbb{R}^r$ and $\epsilon \in \mathbb{R}^n$ are the temporary variables to denote the coefficients for the covariates and residuals respectively. For this latent dimension *d*, we then perform the following update for each gene $j \in \{1, \ldots, p\}$,

$$\hat{Z}_{j,\cdot} \leftarrow \hat{Z}_{j,\cdot} + \beta \cdot \hat{Y}_{j,d}$$
, and $\hat{X}_{\cdot,d} \leftarrow \epsilon$.

It can be seen that after performing this update for every latent dimension $d \in \{1, ..., k\}$, \hat{X} is orthogonal to *C* (i.e., $\hat{X}^{\top}C = 0$) even though the predictive power of our factorization (i.e., $\hat{Y}^{\top}\hat{X} + \hat{Z}^{\top}C$) did not change.

Step 2: We next perform a linear transformation on X̂ and Ŷ. Specifically, using the details in [18], let R = Ŷ^TX̂ ∈ ℝ^{p×n}, with an SVD of R = UDV^T. Then,

$$\hat{X} \leftarrow \left(\frac{n}{p}\right)^{1/4} V \sqrt{D}$$
, and $\hat{Y} \leftarrow \left(\frac{p}{n}\right)^{1/4} U \sqrt{D}$.

It can be seen that after performing this update that both $\hat{X}^{\top}\hat{X}/n$ and $\hat{Y}^{\top}\hat{Y}/p$ are diagonal matrices and are equal even though the predictive power of our factorization (i.e., $\hat{Y}^{\top}\hat{X} + \hat{Z}^{\top}C$) did not change. This ensures identifability of \hat{X} and \hat{Y} .

Estimating overdispersion parameter

We estimate the overdispersion parameter $\gamma_1, \ldots, \gamma_p > 0$, one for each gene. Following our model of the covariate-adjusted sequencing depth (3), we estimate the covariate-adjusted sequencing depth of each cell as

$$\hat{\ell}_{ji} = \exp\left((\hat{Z}_{j,-(cc)})^{\top} C_{i,-(cc)}\right), \quad \text{for all } i \in \{1, \dots, n\}, j \in \{1, \dots, p\}.$$
(13)

Likewise, the estimate of the mean parameter is

$$\hat{\mu}_{ji} = \exp\left((\hat{Y}_{j,\cdot})^{\top}\hat{X}_{i,\cdot} + \hat{Z}_{j,(\mathrm{cc})} \cdot C_{i,(\mathrm{cc})}\right).$$

Then, using a plug-in estimate of the maximum likelihood of the model in (1), we can derive that the estimate of γ_i is

$$\hat{\gamma}_{j} = \max_{\gamma} \sum_{i=1}^{n} A_{ji} \log \hat{\ell}_{i} + \hat{\mu}_{ji} \gamma \log \gamma - \log \Gamma(\gamma \hat{\mu}_{ji}) + \log \Gamma(A_{ji} + \gamma \hat{\mu}_{ji}) - (A_{ji} + \gamma \hat{\mu}_{ji}) \log(\hat{\ell}_{i} + \gamma),$$

where $\Gamma(\cdot)$ is the Gamma function. We estimate this via Newton's method, and this step is primarily handled by eSVD2::estimate nuisance in our codebase.

Computing the posterior distribution

To account for possible model misspecification, we estimate the posterior distribution of the mean and variance of each gene j's expression in cell i. This is based on the model

in (1), where we can leverage the fact that λ_{ji} conditioned on A_{ji} follows a Negative Binomial distribution. Specifically, based on the posterior distribution derived in (4), we compute

$$\hat{\mu}_{ji}^{(\text{post})} = \frac{\hat{\mu}_{ji}/\hat{\gamma}_j + A_{ji}}{1/\hat{\gamma}_j + \hat{\ell}_{ji}}, \quad \text{and} \quad \hat{\nu}_{ji}^{(\text{post})} = \frac{\hat{\mu}_{ji}/\hat{\gamma}_j + A_{ji}}{(1/\hat{\gamma}_j + \hat{\ell}_{ji})^2}, \tag{14}$$

where $\hat{\mu}_{ji}^{(\text{post})}$ and $\hat{v}_{ji}^{(\text{post})}$ is the posterior mean and variance of $\lambda_{ji}|A_{ji}$ for gene *j* in cell *i* respectively. Other work have used similar modeling based on the posterior distribution for single-cell RNA-seq data [23], but not for the purposes of DE testing. This step is primarily handled by eSVD2::compute posterior in our codebase.

Computing the test statistic

In this step, we compute the test statistic, accounting that multiple cells originate from a particular individual. This broadly falls into two steps: (1) computing the mean and variance for a gene among an individual, and (2) computing the test statistic among all the case and control individuals.

With the posterior distribution in (14), we can compute the expected posterior mean and variance for a particular gene *j* in individual *s*. First, we aggregate among the cells for each individual. We assume that after averaging among all the cells from an individual, the resulting distribution can be reasonably approximated by a Gaussian with mean $\hat{\mu}_j^{(s)}$ and variance $\hat{\nu}_j^{(s)}$. Hence, let $\mathcal{I}(s) \subset \{1, \ldots, n\}$ denote the set of cells originating from individual *s*. Then, via large-sample average,

$$\hat{\mu}_j^{(s)} = \frac{\sum_{i \in \mathcal{I}(s)} \hat{\mu}_{ji}^{(\text{post})}}{|\mathcal{I}(s)|}, \quad \text{and} \quad \hat{\nu}_j^{(s)} = \frac{\sum_{i \in \mathcal{I}(s)} \hat{\nu}_{ji}^{(\text{post})}}{|\mathcal{I}(s)|},$$

where $\hat{\mu}_{js}$ and $\hat{\nu}_{js}$ is the posterior mean and variance for gene *j* among all the cells in individual *s* respectively.

Next, we aggregate among individuals. Specifically, among all the case individuals, we can think of the expression of gene *j* as a Gaussian mixture among the individuals. Let A denote the set of case individuals. We then summarize the Gaussian mixture with one Gaussian,

$$\hat{\mu}_j^{(\text{case})} = \frac{\sum_{s \in \mathcal{A}} \hat{\mu}_j^{(s)}}{|\mathcal{A}|}, \quad \text{and} \quad \hat{\nu}_j^{(\text{case})} = \frac{\sum_{s \in \mathcal{A}} \hat{\nu}_j^{(s)}}{|\mathcal{A}|} + \frac{\sum_{s \in \mathcal{A}} (\hat{\mu}_j^{(s)})^2}{|\mathcal{A}|} - \Big(\frac{\sum_{s \in \mathcal{A}} \hat{\mu}_j^{(s)}}{|\mathcal{A}|}\Big)^2.$$

We do an analogous calculation among controls, letting \mathcal{B} denote all the control individuals.

Lastly, we do a two-sample T-test with unequal variances, treating the distribution among cases as a Gaussian with mean $\hat{\mu}_j^{(\text{case})}$ and variance $\hat{\nu}_j^{(\text{case})}$ and the distribution among controls as an analogous Gaussian. Our test statistic for gene *j* is

$$\hat{T}_{j} = \frac{\hat{\mu}_{j}^{(\text{case})} - \hat{\mu}_{j}^{(\text{control})}}{\sqrt{\frac{1}{|\mathcal{A}|} \cdot \hat{\nu}_{j}^{(\text{case})} + \frac{1}{|\mathcal{B}|} \cdot \hat{\nu}_{j}^{(\text{control})}}}.$$
(15)

In our codebase, this step is primarily handled by <code>eSVD2::compute_test_statis-tic</code>.

Performing multiple testing correction

We use the multiple testing procedure based on the empirical null developed in [57] to handle the multiple testing corrections. This framework is also used in other methods designed to test for differential expression from scRNA-seq data (albeit not from cohort data), such as iDEA [58] and SifiNet [59]. We know that biologically, most genes are not directly related to the disease or disorder at the human population level; hence, most genes should be deemed insignificant. The empirical null, where the appropriate null distribution is learned from the data, is a suitable framework to ensure this behavior.

Briefly, we first estimate the degree-of-freedom of each gene $j \in \{1, ..., p\}$ (used for T-tests for unequal variances) by

$$\hat{\mathrm{df}}_{j} = \frac{\left(\frac{1}{|\mathcal{A}|} \cdot \hat{v}_{j}^{(\mathrm{case})} + \frac{1}{|\mathcal{B}|} \cdot \hat{v}_{j}^{(\mathrm{control})}\right)^{2}}{\left(\frac{1}{|\mathcal{A}|} \cdot \hat{v}_{j}^{(\mathrm{case})}\right)^{2} / (|\mathcal{A}| - 1) + \left(\frac{1}{|\mathcal{B}|} \cdot \hat{v}_{j}^{(\mathrm{control})}\right)^{2} / (|\mathcal{B}| - 1)},$$

and convert all the test statistics \hat{T}_i 's into a Z-score via,

$$\hat{Z}_j = \Phi^{-1} \Big(F_{\hat{\mathrm{df}}_j} \big(\hat{T}_j \big) \Big),$$

where $F_b(a)$ is the CDF of the a *t*-distribution with degree freedom b > 0 evaluated at *a*, and $\Phi^{-1}(a)$ is the quantile function of a standard Gaussian evaluated at $a \in (0, 1)$.

To identify the DE genes, we use the locfdr::locfdr function to estimate the empirical null distribution's mean and standard deviation. We then compute the *p*-value of each gene based on the (two-sided) tail area of a Gaussian distribution with that empirical null's mean and standard deviation. We can compute the negative log_{10} *p*-value for visualization purposes in a volcano plot (as in Figs. 3 and 4). To perform multiple testing corrections, we apply stats::p.adjust on these *p*-values via Benjamini-Hochberg. All the genes with an empirical FDR of less than 0.05 are then deemed to be DE genes.

This step is collectively handled by eSVD2::compute_df, eSVD2::compute_test statistic, and eSVD2::compute pvalue in our codebase.

Information about data preprocessing

We describe the high-level ideas on what is needed for our preprocessing of the data prior to using the eSVD-DE. Importantly, we need to: (1) select the cells specific for our analysis, and (2) select the genes of interest for our analysis. For all the datasets in our analysis, the cell type labels were already provided by the author, and we used Seurat::FindVariableFeatures to select most of the genes in our analysis. (We also included housekeeping genes and genes previously reported to be differentially expressed by the authors as well.) No preprocessing of the count data is needed as eSVD-DE models the raw count data directly, using the observed sequencing depth and covariates of the individuals. We more thoroughly detail these steps in the Appendix (Section E) and include details on which covariates were included in all our analyses as well as the parameters used for eSVD-DE.

Supplementary Information

The online version contains supplementary material available at https://doi.org/10.1186/s12859-024-05724-7.

Additional file 1. High-level description on what is statistically challenging about doing differential expression testing after denoising a datset, simulation details and results, and supplemental details of the lung, colon, and brain analyses, and supplemental figures A1-F12.

Additional file 2. This file contains the analysis of T-cells from either the Adams data (from [6]) or the Habermann data (from [19]). It contains the estimated log-fold change after removing confounding effects, the eSVD-DE test statistic, the negative log10 p-value, and Benjamini-Hochberg adjusted p-value, and whether the gene is deemed to be a DE gene.

Additional file 3. This file contains the analysis of cycling TA, enterocyte progenitors, TA 1, and TA 2 when analyzing either inflamed or non-inflamed tissues from the Smillie dataset (from [20]). It contains the estimated log-fold change after removing confounding effects, the eSVD-DE test statistic, the negative log10 p-value, and Benjamini-Hochberg adjusted p-value, and whether the gene is deemed to be a DE gene.

Additional file 4. This file contains the analysis of astrocytes, endothelial cells, IN-SST, IN-VIP, layer 2/3, layer 4, layer 5/6, layer 5/6-CC, microglia, oligodendrocytes, and OPC from the Velmeshev dataset (from [21]). It contains the estimated log-fold change after removing confounding effects, the eSVD-DE test statistic, the negative log10 p-value, and Benjamini-Hochberg adjusted p-value, and whether the gene is deemed to be a DE gene.

Acknowledgements

We thank Timothy Barry, Jing Lei, Ronald Yurko, and Nancy Zhang for useful comments when developing this method.

Author Contributions

KZL and KR conceived of the idea. KZL and YQ coded eSVD-DE in R and C++. KZL performed the analyses. KZL, YQ and KR wrote the paper.

Funding

This project was funded by National Institute of Mental Health (NIMH) Grant R01MH123184.

Availability of data and materials

All analyses were done on publicly available data. The Adams data (from [6]) is downloaded from GSE136831, and the individual covariates were provided by the authors of the paper. The Habermann data (from [19]) is downloaded from GSE135893, and the individual covariates were provided by the authors of the paper. The Smilie data (from [20]) is downloaded froms https://github.com/cssmillie/ulcerative_colitis, while the individuals' covariates and DE genes are downloaded as Table S1 and S4 from https://www. ncbi.nlm.nih.gov/pmc/articles/PMC6662628/. The Velmeshev data (from [21]) is downloaded from https://cells.ucsc.edu/?ds=autism, while the DE genes is downloaded as Data S4 from http://science.sciencemag.org/content/suppl/2019/05/15/364.6441.685.DC1. The housekeeping genes (from [40]) is downloaded as Supplementary Table 1 from https:// academic.oup.com/nar/article/49/D1/D947/5871367#supplementary-data. The SFARI genes were downloaded from https://gene.sfari.org/database/gene-scoring/ on January 6, 2022 (using the September 2, 2021 release). The Gandal genes (from [49]) is downloaded as Supplementary Data 3 from https://www.nature.com/articles/s41586-022-05377-7, where we look at the DEGene Statistics sheet and select the genes with WholeCortex ASD FDR less than 0.05.Code availability See https://github.com/linnykos/eSVD2 for the R package eSVD2 that contains all the functions used in for the eSVD-DE analysis. See https://github.com/linnykos/eSVD2_examples for all the scripts used in the analysis.

Declarations

Ethics approval Not applicable.

Consent to participate Not applicable.

Consent for publication Not applicable.

Competing interests

The authors declare that they have no competing interests.

Received: 6 December 2023 Accepted: 28 February 2024 Published online: 15 March 2024

References

- 1. Kim N, Kim HK, Lee K, Hong Y, Cho JH, Choi JW, Lee J-I, Suh Y-L, Ku BM, Eum HH. Single-cell RNA sequencing demonstrates the molecular and cellular reprogramming of metastatic lung adenocarcinoma. Nat Commun. 2020;11(1):1–15.
- Gabitto M, Travaglini K, Ariza J, Kaplan E, Long B, Rachleff V, Ding Y, Mahoney J. Dee N, Goldy J, others Haynor D, Gatto NM, Jayadev S, Mutfi S, Ng L, Mukherjee S, Crane PK, Latimer CS, Levi BP, Smith K, Close JL, Miller JA, Hodge RD, Larson EB, Grabowski TJ, Hawrylycz M, Keene CD, Lein ES. Integrated multimodal cell atlas of Alzheimer disease 2023.
- Perez RK, Gordon MG, Subramaniam M, Kim MC, Hartoularos GC, Targ S, Sun Y, Ogorodnikov A, Bueno R, Lu A. Single-cell RNA-seg reveals cell type-specific molecular and genetic associations to lupus. Science. 2022;376(6589):1970.
- Squair JW, Gautier M, Kathe C, Anderson MA, James ND, Hutson TH, Hudelle R, Qaiser T, Matson KJ, Barraud Q, Barraud Q, Levine AJ, La Manno G, Skinnider MA, Courtine G. Confronting false discoveries in single-cell differential expression. Nat Commun. 2021;12(1):5692.
- Mallick H, Chatterjee S, Chowdhury S, Chatterjee S, Rahnavard A, Hicks SC. Differential expression of single-cell RNA-seq data using tweedie models. Stat Med. 2022;41(18):3492–510.
- Adams TS, Schupp JC, Poli S, Ayaub EA, Neumark N, Ahangari F, Chu SG, Raby BA, Deluliis G, Januszyk M, Duan Q, Arnett HA, Siddiqui A, Washko GR, Homer R, Yan X, Rosas IO, Kaminski N. Single-cell RNA-seq reveals ectopic and aberrant lungresident cell populations in idiopathic pulmonary fibrosis. Sci Adv. 2020;6(28):1983.
- Auerbach BJ, Hu J, Reilly MP, Li M. Applications of single-cell genomics and computational strategies to study common disease and population-level variation. Genome Res. 2021;31(10):1728–41.
- Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. Genome Biol. 2014;15(12):550.
- McCarthy DJ, Chen Y, Smyth GK. Differential expression analysis of multifactor RNA-seq experiments with respect to biological variation. Nucl Acids Res. 2012;40(10):4288–97.
- Finak G, McDavid A, Yajima M, Deng J, Gersuk V, Shalek AK, Slichter CK, Miller HW, McElrath MJ, Prlic M, Linsley PS, Gottardo R. MAST: A flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data. Genome Biol. 2015;16(1):278.
- 11. He L, Davila-Velderrain J, Sumida TS, Hafler DA, Kellis M, Kulminski AM. NEBULA is a fast negative binomial mixed model for differential or co-expression analysis of large-scale multi-subject single-cell data. Commun Biol. 2021;4(1):629.
- Townes FW, Hicks SC, Aryee MJ, Irizarry RA. Feature selection and dimension reduction for single-cell RNA-seq based on a multinomial model. Genome Biol. 2019;20(1):1–16.
- Risso D, Perraudeau F, Gribkova S, Dudoit S, Vert J-P. A general and flexible method for signal extraction from single-cell RNA-seg data. Nat Commun. 2018;9(1):284.
- 14. Nicol PB, Miller JW. Model-based dimensionality reduction for single-cell RNA-seq using generalized bilinear models. bioRxiv. 2023;2023–04
- 15. And rews TS, Hemberg M. False signals induced by single-cell imputation. F1000Research. 2018;7:1.
- Boyeau P, Regier J, Gayoso A, Jordan MI, Lopez R, Yosef N. An empirical bayes method for differential expression analysis of single cells with deep generative models. Proc Natl Acad Sci. 2023;120(21):2209124120.
- 17. Weinberger E, Lin C, Lee SI. Isolating salient variations of interest in single-cell data with contrastiveVI. Nat Methods. 2023;1–10
- Lin KZ, Lei J, Roeder K. Exponential-family embedding with application to cell developmental trajectories for single-cell RNA-seq data. J Am Stat Assoc. 2021;116(534):457–70.
- Habermann AC, Gutierrez AJ, Bui LT, Yahn SL, Winters NI, Calvi CL, Peter L, Chung M-I, Taylor CJ, Jetter C, Raju L, Roberson J, Ding G, Wood L, Sucre JMS, Richmond BW, Serezani AP, McDonnell WJ, Mallal SB, Bacchetta MJ, Loyd JE, Shaver CM, Ware LB, Bremner R, Walia R, Blackwell TS, Banovich NE, Kropski JA. Single-cell RNA sequencing reveals profibrotic roles of distinct epithelial and mesenchymal lineages in pulmonary fibrosis. Sci Adv. 2020;6(28):1972.
- Smillie CS, Biton M, Ordovas-Montanes J, Sullivan KM, Burgin G, Graham DB, Herbst RH, Rogel N, Slyper M, Waldman J, Sud M, Andrews E, Velonias G, Haber AL, Jagadeesh K, Vickovic S, Yao J, Stevens C, Dionne D, Nguyen LT, Villani A-C, Hofree M, Creasey EA, Huang H, Rozenblatt-Rosen O, Garber JJ, Khalili H, Desch AN, Daly MJ, Ananthakrishnan AN, Shalek AK, Xavier RJ, Regev A. Intra-and inter-cellular rewiring of the human colon during ulcerative colitis. Cell. 2019;178(3):714–30.
- Velmeshev D, Schirmer L, Jung D, Haeussler M, Perez Y, Mayer S, Bhaduri A, Goyal N, Rowitch DH, Kriegstein AR. Singlecell genomics identifies cell type-specific molecular changes in autism. Science. 2019;364(6441):685–9.

- Sarkar A, Stephens M. Separating measurement and expression models clarifies confusion in single cell RNA-seq analysis. Nat Genet. 2021;53(6):770–7.
- Huang M, Wang J, Torre E, Dueck H, Shaffer S, Bonasio R, Murray JI, Raj A, Li M, Zhang NR. SAVER: gene expression recovery for single-cell RNA sequencing. Nat Methods. 2018;15(7):539–42.
- 24. Gayoso A, Steier Z, Lopez R, Regier J, Nazor KL, Streets A, Yosef N. Joint probabilistic modeling of single-cell multi-omic data with totalVI. Nat Methods. 2021;18(3):272–82.
- Iterson M, Zwet EW, Heijmans BT. Controlling bias and inflation in epigenome-and transcriptome-wide association studies using the empirical null distribution. Genome Biol. 2017;18(1):1–13.
- Chen W, Li Y, Easton J, Finkelstein D, Wu G, Chen X. UMI-count modeling and differential expression analysis for singlecell RNA sequencing. Genome Biol. 2018;19(1):70.
- Zhang M, Liu S, Miao Z, Han F, Gottardo R, Sun W. IDEAS: Individual level differential expression analysis for single-cell RNA-seq data. Genome Biol. 2022;23(1):1–17.
- Korthauer KD, Chu L-F, Newton MA, Li Y, Thomson J, Stewart R, Kendziorski C. A statistical approach for identifying differential distributions in single-cell RNA-seq experiments. Genome Biol. 2016;17(1):1–15.
- Schefzik R, Flesch J, Goncalves A. Fast identification of differential distributions in single-cell RNA-sequencing data with waddR. Bioinformatics. 2021;37(19):3204–11.
- Junttila S, Smolander J, Elo LL. Benchmarking methods for detecting differential states between conditions from multisubject single-cell RNA-seq data. Brief Bioinf. 2022;23(5):286.
- 31. Liu Y, Zhao J, Adams TS, Wang N, Schupp JC, Wu W, McDonough JE, Chupp GL, Kaminski N, Wang Z, Yan X. iDESC: identifying differential expression in single-cell RNA sequencing data with multiple subjects. BMC Bioinf. 2023;24(1):318.
- Tenenbaum JB, Silva VD, Langford JC. A global geometric framework for nonlinear dimensionality reduction. Science. 2000;290(5500):2319–23.
- Kraemer G, Reichstein M, Mahecha MD. dimRed and coRanking unifying dimensionality reduction in R. R J. 2018;10(1):342–58.
- 34. McInnes L, Healy J, Melville J. UMAP: Uniform manifold approximation and projection for dimension reduction. arXiv preprint arXiv:1802.03426 (2018)
- 35. Ham J, Lee DD, Mika S, Schölkopf B A kernel view of the dimensionality reduction of manifolds. In: Proceedings of the Twenty-first International Conference on Machine Learning. 2004; p. 47
- Wu H-T, Wu N. Think globally, fit locally under the manifold setup: asymptotic analysis of locally linear embedding. Ann Stat. 2018;46(6B):3805–37.
- 37. Perturbation bounds for procrustes, classical scaling, and trilateration, with applications to manifold learning. J Machine Learn Res. 2020;21.
- Li Y, Ge X, Peng F, Li W, Li JJ. Exaggerated false positives by popular differential expression methods when analyzing human population samples. Genome Biol. 2022;23(1):1–13.
- 39. Hafemeister C, Satija R. Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. Genome Biol. 2019;20(1):1–15.
- Hounkpe BW, Chenou F, Lima F, De Paula EV. HRT atlas v1.0 database Redefining human and mouse housekeeping genes and candidate reference transcripts by mining massive RNA-seq datasets. Nucl Acids Res. 2021;49(D1):947–55.
- 41. Risso D, Ngai J, Speed TP, Dudoit S. Normalization of RNA-seq data using factor analysis of control genes or samples. Nat Biotechnol. 2014;32(9):896–902.
- 42. Lause J, Berens P, Kobak D. Analytic Pearson residuals for normalization of single-cell RNA-seq UMI data. Genome Biol. 2021;22(1):1–20.
- Cole MB, Risso D, Wagner A, DeTomaso D, Ngai J, Purdom E, Dudoit S, Yosef N. Performance assessment and selection of normalization procedures for single-cell RNA-seq. Cell Syst. 2019;8(4):315–28.
- Lun AT, Bach K, Marioni JC. Pooling across cells to normalize single-cell RNA sequencing data with many zero counts. Genome Biol. 2016;17(1):75.
- Nance T, Smith KS, Anaya V, Richardson R, Ho L, Pala M, Mostafavi S, Battle A, Feghali-Bostwick C, Rosen G, Montgomery SB. Transcriptome analysis reveals differential splicing events in IPF lung tissue. PLoS ONE. 2014;9(3):92111.
- 46. Joshi N, Watanabe S, Verma R, Jablonski RP, Chen CI, Cheresh P, Markov NS, Reyfman PA, McQuattie-Pimentel, AC, Sichizya L, Lu Z, Piseaux R, Kirchenbuechler D, Flozak AS, Gottardi CJ, Cuda CM, Perlman H, Jain M, Kamp DW, Budinger GRS, Misharin AV. A spatially restricted fibrotic niche in pulmonary fibrosis is sustained by M-CSF/M-CSFR signalling in monocyte-derived alveolar macrophages. Eur Respirat J. 2020;55(1)
- 47. Gauldie J, Kolb M, Ask K, Martin G, Bonniaud P, Warburton D. Smad3 signaling involved in pulmonary fibrosis and emphysema. Proc Am Thorac Soc. 2006;3(8):696–702.
- 48. SFARI gene database (2022). https://gene.sfari.org/. Accessed 2022 October 20.
- Gandal MJ, Haney JR, Wamsley B, Yap CX, Parhami S, Emani PS, Chang N, Chen GT, Hoftman GD, Alba D, Ramaswami G, Hartl CL, Bhattacharya A, Luo C, Jin T, Wang D, Kawaguchi R, Quintero D, Ou J, Wu YE, Parikshak NN, Swarup V, Belgard TG, Gerstein M, Pasaniuc B, Geschwind DH. Broad transcriptomic dysregulation occurs across the cerebral cortex in ASD. Nature. 2022;611(7936):532–9. https://doi.org/10.1038/s41586-022-05377-7.
- 50. Efron B. Microarrays, empirical Bayes and the two-groups model. Stat Sci. 2008;23(1):1–22.
- Zhao B-W, Su X-R, Hu P-W, Ma Y-P, Zhou X, Hu L. A geometric deep learning framework for drug repositioning over heterogeneous information networks. Brief Bioinf. 2022;23(6):384.
- 52. Zhao B-W, Su X-R, Hu P-W, Huang Y-A, You Z-H, Hu L. iGRLDTI: an improved graph representation learning method for predicting drug-target interactions over heterogeneous biological information network. Bioinformatics. 2023;39(8):451.
- 53. Read DF, Daza, RM, Booth GT, Jackson DL, Gladden RG, Srivatsan SR. Ewing B, Franks JM, Spurrell CH. Gomes AR, O'Day D, Gogate AA, Martin BK, Starita L, Lin Y, Shendure J, Lin S, Trapnell C Single-cell analysis of chromatin and expression reveals age-and sex-associated alterations in the human heart. bioRxiv; 2022.
- 54. Agarwal D, Wang J, Zhang NR. Data denoising and post-denoising corrections in single cell RNA sequencing. Stat Sci. 2020;35(1):112–28.
- 55. Wang L, Zhang X, Gu Q. A unified computational and statistical framework for nonconvex low-rank matrix estimation. Artif Intell Stat. 2017;981–990. PMLR.

- Li X, Lu J, Arora R, Haupt J, Liu H, Wang Z, Zhao T. Symmetry, saddle points, and global optimization landscape of nonconvex matrix factorization. IEEE Trans Inf Theory. 2019;65(6):3489–514.
- 57. Efron B. Large-scale simultaneous hypothesis testing: the choice of a null hypothesis. J Am Stat Assoc. 2004;99(465):96–104.
- Ma Y, Sun S, Shang X, Keller ET, Chen M, Zhou X. Integrative differential expression and gene set enrichment analysis using summary statistics for scRNA-seq studies. Nat Commun. 2020;11(1):1–13.
- Gao Q, Ji, Z, Wang L, Owzar, K, Li QJ, Chan C, Xie J. SifiNet: a robust and accurate method to identify feature gene sets and annotate cells. bioRxiv, 2023;2023–05

Publisher's Note

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