

# Feature Selection and Dimension Reduction for Single Cell RNA-Seq based on a Multinomial Model

F. WILLIAM TOWNES<sup>1</sup>, STEPHANIE C. HICKS<sup>2</sup>,  
MARTIN J. ARYEE<sup>1,3,4,5</sup>, RAFAEL A. IRIZARRY<sup>\*,1,6</sup>

<sup>1</sup>Department of Biostatistics, Harvard University, Boston, MA,

<sup>2</sup>Department of Biostatistics, Johns Hopkins University, Baltimore, MD,

<sup>3</sup>Molecular Pathology Unit, Massachusetts General Hospital, Charlestown, MA,

<sup>4</sup>Center for Cancer Research, Massachusetts General Hospital, Charlestown, MA,

<sup>5</sup>Department of Pathology, Harvard Medical School, Boston, MA,

<sup>6</sup>Department of Data Sciences, Dana-Farber Cancer Institute, Boston, MA

ftownes@g.harvard.edu, shicks19@jhu.edu, arye.martin@mgh.harvard.edu,

rafa@jimmy.harvard.edu

March 11, 2019

---

\* To whom correspondence should be addressed.

### Abstract

Single cell RNA-Seq (scRNA-Seq) profiles gene expression of individual cells. Recent scRNA-Seq datasets have incorporated unique molecular identifiers (UMIs). Using negative controls, we show UMI counts follow multinomial sampling with no zero-inflation. Current normalization procedures such as log of counts per million and feature selection by highly variable genes produce false variability in dimension reduction. We propose simple multinomial methods, including generalized principal component analysis (GLM-PCA) for non-normal distributions, and feature selection using deviance. These methods outperform current practice in a downstream clustering assessment using ground-truth datasets.

Keywords: gene expression, single cell, RNA-Seq, dimension reduction, variable genes, principal components analysis, GLM-PCA

## 1 Background

Single cell RNA-Seq (scRNA-Seq) is a powerful tool for profiling gene expression patterns in individual cells, facilitating a variety of analyses such as identification of novel cell types [1, 2]. In a typical protocol, single cells are isolated in liquid droplets and messenger RNA (mRNA) is captured from each cell, converted to cDNA by reverse transcriptase (RT), then amplified using polymerase chain reaction (PCR) [3, 4, 5]. Finally, fragments are sequenced and expression of a gene in a cell is quantified by the number of sequencing reads that mapped to that gene [6]. A crucial difference between scRNA-Seq and traditional bulk RNA-Seq is the low quantity of mRNA isolated from individual cells, which requires a larger number of PCR cycles to produce enough material for sequencing (bulk RNA-Seq comingles thousands of cells per sample). Thus, many of the reads counted in scRNA-Seq are duplicates of a single mRNA molecule in the original cell [7]. Early scRNA-Seq studies using protocols such as SMART-Seq2 [8] analyzed these *read counts* directly, and several methods were developed to facilitate this [9]. However, newer protocols typically include unique molecular identifiers (UMIs) which enable computational removal of PCR duplicates [10], producing *UMI counts*. Although a zero UMI count is equivalent to a zero read count, nonzero read counts are larger than their corresponding UMI counts. In general, all scRNA-Seq data contain large numbers of zero counts (often  $> 90\%$  of the data, sometimes called *dropouts*). Here, we focus on the analysis of scRNA-Seq data with UMI counts.

Starting from raw counts, a scRNA-Seq data analysis typically includes normalization, feature selection, and dimension reduction steps. Normalization seeks to adjust for differences in experimental conditions between samples (individual cells), so that these do not confound true biological differences. For example, the efficiency of mRNA capture and RT is variable between samples (technical variation), causing different cells to have different total UMI counts, even if the number of molecules in the original cells is identical. Feature selection refers to excluding uninformative genes such as those which exhibit no meaningful biological variation across samples. Since scRNA-Seq experiments usually

examine cells within a single tissue, only a small fraction of genes are expected to be informative since many genes are biologically variable only across different tissues. Dimension reduction aims to embed each cell's high-dimensional expression profile into a low-dimensional representation to facilitate visualization and clustering.

While a plethora of methods [11, 12, 13, 5, 14] have been developed for each of these steps, here we describe what is considered to be the standard pipeline [14]. First, raw counts are normalized by scaling of sample-specific *size factors*, followed by log-transformation, which attempts to reduce skewness. Next, feature selection involves identifying the top 500-2,000 genes by computing either their coefficient of variation (highly variable genes [15, 16]), or average expression level (highly expressed genes) across all cells [14]. Alternatively, highly dropout genes may be retained [17]. Principal component analysis (PCA) [18] is the most popular dimension reduction method (see for example tutorials for Seurat [16] and Cell Ranger [5]). PCA compresses each cell's 2,000-dimensional expression profile into, say, a 10-dimensional vector of principal component coordinates or latent factors. Prior to PCA, data are usually centered and scaled so that each gene has mean zero and standard deviation one (*z-score* transformation). Finally, a clustering algorithm can be applied to group cells with similar representations in the low-dimensional PCA space.

Despite the appealing simplicity of this standard pipeline, the characteristics of scRNA-Seq UMI counts present difficulties at each stage. Many normalization schemes derived from bulk RNA-Seq cannot compute size factors stably in the presence of large numbers of zeros [19]. A numerically stable and popular method is to set the size factor for each cell as the total counts divided by  $10^6$  (*counts per million*, CPM). Note that CPM does not alter zeros, which dominate scRNA-Seq data. Log-transformation is not possible for exact zeros so it is common practice to add a small *pseudocount* such as one to all normalized counts prior to taking the log. The choice of pseudocount is arbitrary and can introduce subtle biases in the transformed data [20]. For a statistical interpretation of the pseudocount see Section 4.2. Similarly, the use of highly variable genes for feature selection is somewhat arbitrary since the observed variability will depend on the pseudocount: pseudocounts close to zero arbitrarily increase the variance of genes with zero counts. Finally, PCA implicitly relies on Euclidean geometry, which may not be appropriate for highly sparse, discrete, and skewed data, even after normalizations and transformations [21].

Widely used methods for analysis of scRNA-Seq lack statistically rigorous justification based on a plausible data generating mechanism for UMI counts. Instead, it appears many of the techniques have been borrowed from the data analysis pipelines developed for read counts, especially those based on bulk RNA-Seq [22]. For example, models based on the lognormal distribution cannot account for exact zeros, motivating the development of zero-inflated lognormal models for scRNA-Seq read counts [23, 24, 25, 26]. Alternatively, ZINB-WAVE uses a zero-inflated negative binomial model for dimension reduction of read counts [27]. However, as shown below, the sampling distribution of UMI counts differs markedly from read counts, so application of read count models to UMI

counts needs either theoretical or empirical justification.

We present a unifying statistical foundation for scRNA-Seq with UMI counts based on the multinomial distribution. The multinomial model adequately describes negative control data and there is no need to model zero inflation. We show the mechanism by which PCA on log-normalized UMI counts can lead to distorted low-dimensional factors and false discoveries. We identify the source of the frequently observed and undesirable fact that the fraction of zeros reported in each cell drives the first principal component in most experiments [28]. To remove these distortions, we propose the use of GLM-PCA, a generalization of PCA to exponential family likelihoods [29]. GLM-PCA operates on raw counts, avoiding the pitfalls of normalization. We also demonstrate that applying PCA to deviance or Pearson residuals provides a useful and fast approximation to GLM-PCA. We provide a closed-form deviance statistic as a feature selection method. We systematically compare the performance of all combinations of methods using three ground-truth datasets and assessment procedures from [14]. We conclude by suggesting best practices.

## 2 Results and discussion

### 2.1 Datasets

We used five public 10x genomics UMI counts datasets from [5] to benchmark our methods. The first dataset is a highly controlled experiment specifically designed to understand technical variability. No actual cells were used to generate this dataset. Instead, each of the 1,015 droplets received the same ratio of 92 synthetic spike-in RNA molecules from External RNA Controls Consortium (ERCC). We refer to this dataset as the *technical replicates negative control* as there is no biological variability whatsoever and, in principle, each expression profile should be the same.

The second dataset was generated by processing a homogeneous population of 2,612 monocyte cells. The cells were purified using fluorescence activated cell sorting (FACS). We refer to this dataset as the *biological replicates negative control*. Because these cells were all the same type, we did not expect to observe any significant differences in unsupervised analysis.

The third dataset consists of 68,579 fresh, unsorted peripheral blood mononuclear cells (PBMCs). This dataset was used to compare computational speed of different dimension reduction algorithms. We refer to it as the PBMC 68K dataset.

The remaining two datasets were created by [14]. In the Zheng 4eq dataset, there are 3,994 PBMCs divided equally into four cell types. In the Zheng 8eq dataset, there are 3,994 PBMCs divided equally into eight cell types. In these positive control datasets, the cluster identity of all cells was assigned independently of gene expression (using FACS), so they served as ground truth labels.

## 2.2 UMI count distribution differs from reads

To illustrate the marked difference between UMI count distributions and read count distributions, we created histograms from individual genes and spike-ins of the negative control data. Here, the UMI counts are the computationally de-duplicated versions of the read counts; both measurements are from the same experiment, so no differences are due to technical or biological variation. The results suggest that while read counts appear zero-inflated and multimodal, UMI counts follow a discrete distribution with no zero inflation (Figure S1). The apparent zero inflation in read counts is a result of PCR duplicates.

## 2.3 Multinomial sampling distribution for UMI counts

Consider a single cell  $i$  containing  $t_i$  total mRNA transcripts. Let  $n_i$  be the total number of UMIs for the same cell. When the cell is processed by a scRNA-Seq protocol, it is lysed, then some fraction of the transcripts are captured by beads within the droplets. A series of complex biochemical reactions occur, including attachment of barcodes and UMIs, and reverse transcription of the captured mRNA to a cDNA molecule. Finally, the cDNA is sequenced and PCR duplicates are removed to generate the UMI counts [5]. In each of these stages, some fraction of the molecules from the previous stage are lost [30, 5, 7]. In particular, reverse transcriptase is an inefficient and error-prone enzyme [31]. Therefore the number of UMI counts representing the cell is much less than the number of transcripts in the original cell ( $n_i \ll t_i$ ). Specifically,  $n_i$  typically ranges from 1,000 – 10,000 while  $t_i$  is estimated to be approximately 200,000 for a typical mammalian cell [32]. Furthermore, which molecules are selected and successfully become UMIs is a random process. Let  $x_{ij}$  be the true number of mRNA transcripts of gene  $j$  in cell  $i$ , and  $y_{ij}$  be the UMI count for the same gene and cell. We define the *relative abundance*  $\pi_{ij}$  as the true number of mRNA transcripts represented by gene  $j$  in cell  $i$  divided by the total number of mRNA transcripts in cell  $i$ . Relative abundance is given by  $\pi_{ij} = x_{ij}/t_i$  where total transcripts  $t_i = \sum_j x_{ij}$ . Since  $n_i \ll t_i$ , there is a “competition to be counted” [33]; genes with large relative abundance  $\pi_{ij}$  in the original cell are more likely to have nonzero UMI counts, but genes with small relative abundances may be observed with UMI counts of exact zeros. The UMI counts  $y_{ij}$  are a multinomial sample of the true biological counts  $x_{ij}$ , containing only relative information about expression patterns in the cell [34, 33].

The multinomial distribution can be approximated by independent Poisson distributions, and overdispersed multinomials by independent negative binomial distributions. These approximations are useful for computational tractability. Details are provided in the Methods.

The multinomial model makes two predictions which we verified using negative control data. First, the fraction of zeros in a sample (cell or droplet) is inversely related to the total number of UMIs in that sample. Second, the probability of an endogenous gene or ERCC spike-in having zero counts is a decreasing function of its mean expression (equations provided in Methods). Both

of these predictions were validated by the negative control data (Figure 1). In particular, the empirical probability of a gene being zero across droplets was well calibrated to the theoretical prediction based on the multinomial model. This also demonstrates that UMI counts are not zero inflated.

These results are consistent with [35], which also found that the relationship between average expression and zero probability follows the theoretical curve predicted by a Poisson model using negative control data processed with Indrop [4] and Dropseq [3] protocols. This suggests our data generating mechanism is an accurate model of technical noise in real data.

## 2.4 Normalization and log transformation distorts UMI data

Standard scRNA-Seq analysis involves normalizing raw counts using size factors, applying a log transformation with a pseudocount, and then centering and scaling each gene before dimension reduction. The most popular normalization is counts per million (CPM). The CPM are defined as  $(y_{ij}/n_i) * 10^6$  (i.e. the size factor is  $n_i/10^6$ ). This is equivalent to the MLE for relative abundance  $\hat{\pi}_{ij}$  multiplied by  $10^6$ . The log-CPM are then  $\log_2(c + \hat{\pi}_{ij}10^6) = \log_2(\hat{\pi}_{ij}) + C$ , where  $\hat{\pi}_{ij}$  is a maximum a posteriori estimator (MAP) for  $\pi_{ij}$  (mathematical justification and interpretation of this approach provided in Methods). The additive constant  $C$  is irrelevant if data are centered for each gene after log transformation, as is common practice. Thus, normalization of raw counts is equivalent to using MLEs or MAP estimators of the relative abundances.

Log transformation of MLEs is not possible for UMI counts due to exact zeros, while log transformation of MAP estimators of  $\pi_{ij}$  systematically distorts differences between zero and nonzero UMI counts, depending on the arbitrary pseudocount  $c$  (derivations provided in Methods). To illustrate this phenomenon, we examined the distribution of an illustrative gene before and after the log transform with varying normalizations using the biological replicates negative control data (Figure 2). Consistent with our theoretical predictions, this artificially caused the distribution to appear zero inflated and exaggerated differences between cells based on whether the count was zero or nonzero.

Focusing on the entire negative control datasets, we applied PCA to the log transformed CPMs and observed a strong correlation between the first principal component (PC) and the fraction of zeros, consistent with [28]. Additionally, the first PC correlates with the log of total UMI, which is consistent with the multinomial model (Figure 3). Based on these results, the log transformation is not necessary and in fact detrimental for analysis of UMI counts. The benefits of avoiding normalization by instead directly modeling raw counts have been demonstrated in the context of differential expression [36]. Where normalization is unavoidable, we propose the use of approximate multinomial deviance residuals (defined in Section 4.4) instead of log-transformed CPM.

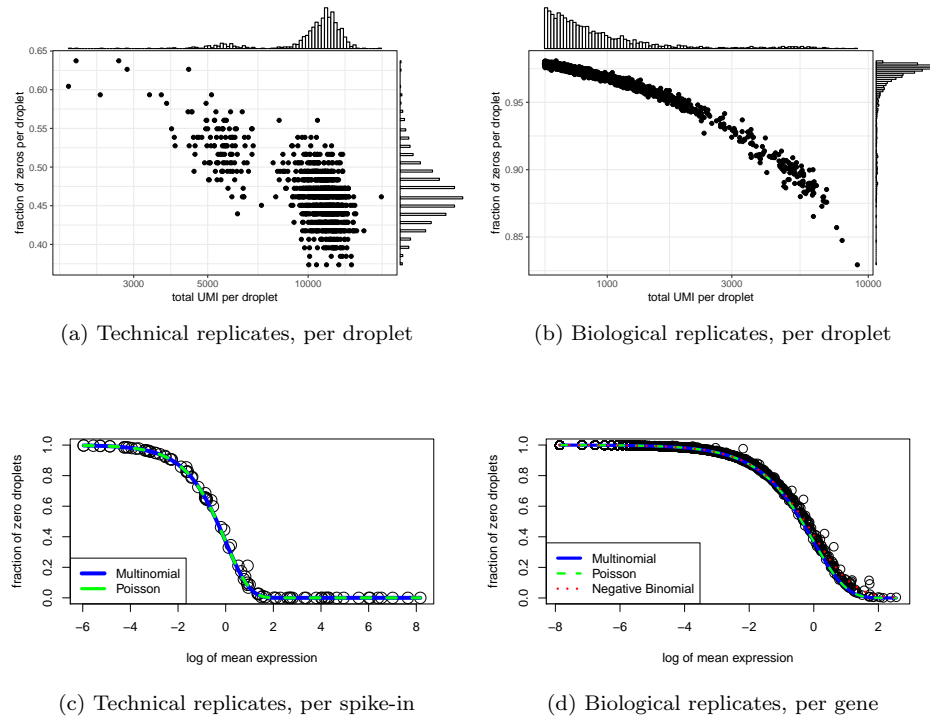


Figure 1: Multinomial model adequately characterizes sampling distributions of technical and biological replicates negative control data. a) Fraction of zeros is plotted against the total number of UMI in each droplet for the technical replicates. b) As a) but for cells in the biological replicates. c) After down-sampling replicates to 10,000 UMIs per droplet to remove variability due to differences in sequencing depth, the fraction of zeros is computed for each gene and plotted against the log of expression across all samples for the technical replicates data. The solid curve is theoretical probability of observing a zero as a function of the expected counts derived from the multinomial model (blue) and its Poisson approximation (green). d) As c) but for the biological replicates dataset and after down-sampling to 575 UMIs per cell. Here we also add the theoretical probability derived from a negative binomial model (red).

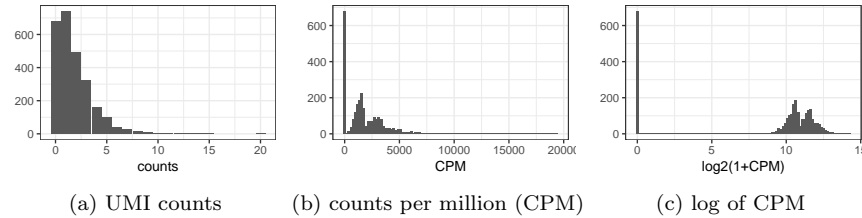


Figure 2: Example of how current approaches to normalization and transformation artificially distort differences between zero and nonzero counts. a) UMI count distribution for gene ENSG00000114391 in the biological replicates negative control dataset. b) Counts per million (CPM) distribution for the exact same count data. c) Distribution of  $\log_2(1 + CPM)$  values for the exact same count data.

## 2.5 Zero inflation is an artifact of log-normalization

To see how normalization and log-transformation introduce the appearance of zero inflation, consider the following example. Let  $y_{ij}$  be the observed UMI counts following a multinomial distribution with size  $n_i$  for each cell and relative abundance  $\pi_j$  for each gene, constant across cells. Focusing on a single gene  $j$ ,  $y_{ij}$  follows a binomial distribution with parameters  $n_i, p_j$ . Assume  $\pi_j = 10^{-4}$  and the  $n_i$  range from 1,000 – 3,000, which is consistent with the biological replicates negative control data (Figures S1 and 1). Under this assumption we expect to see about 74-90% zeros, 22-30% ones, and less than 4% values above one. However, notice that after normalization to CPM and log transformation, all the zeros remain  $\log_2(1 + 0) = 0$ , yet the ones turn into values ranging from  $\log_2(1 + 1/3000 * 10^6) = \log_2(334) \approx 8.4$  to  $\log_2(1001) \approx 10$ . The few values that are 2 will have values ranging from  $\log_2(668) \approx 9.4$  to  $\log_2(2001) \approx 11$ . The large, artificial gap between zero and nonzero values makes the log-normalized data appear zero-inflated (Figure 2). The variability in CPM values across cells is almost completely driven by the variability in  $n_i$ . Indeed, it shows up as the primary source of variation in PCA plots (Figure 3).

## 2.6 Generalized PCA for dimension reduction of sparse counts

While PCA is a popular dimension reduction method, it is implicitly based on Euclidean distance, which corresponds to maximizing a Gaussian likelihood. Since UMI counts are not normally distributed, even when normalized and log transformed, this distance metric is inappropriate, causing PCA to produce distorted latent factors (Figure 3). We propose the use of PCA for generalized linear models (GLMs) [29], or GLM-PCA as a more appropriate alternative. The GLM-PCA framework allows for a wide variety of likelihoods suitable for data types such as counts and binary values. While the multinomial likelihood is ideal for modeling technical variability in scRNA-Seq UMI counts (Figure 1),



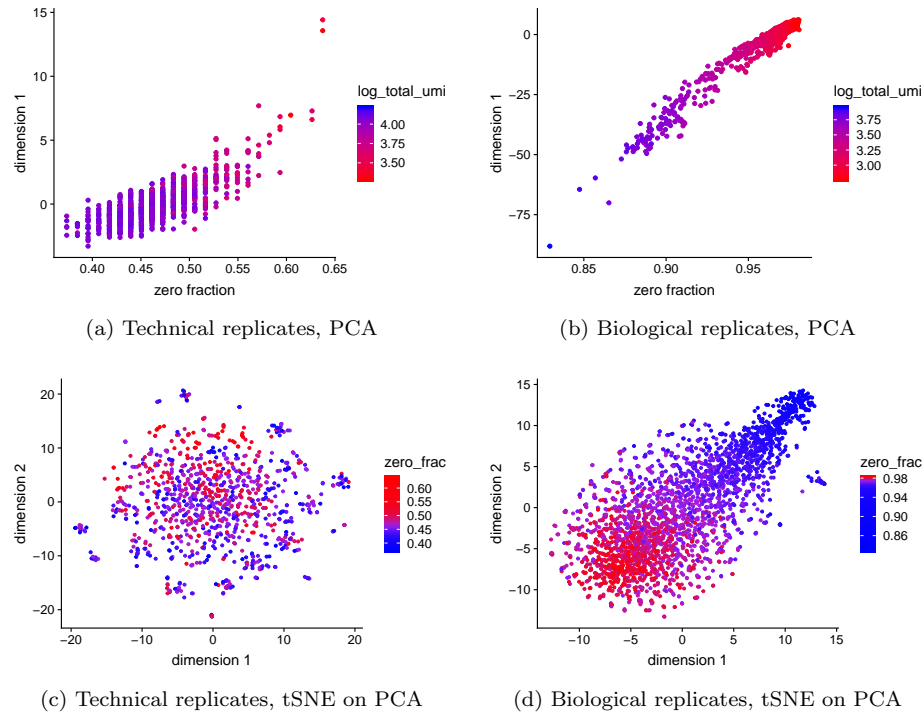


Figure 3: Current approaches to normalization and transformation induce variability in the fraction of zeros across cells to become the largest source of variability which in turn biases clustering algorithms to produce false positive results based on distorted latent factors. a) First principal component (PC) from the technical replicates dataset plotted against fraction of zeros for each cell. A red to blue color scale represents total UMIs per cell. b) as a) but for the biological replicates data. c) Using the technical replicates, we applied t-distributed stochastic neighbor embedding (tSNE) with perplexity 30 to the top 50 PCs computed from log-CPM. The first two tSNE dimensions are shown with a blue to red color scale representing the fraction of zeros. d) as c) but for the biological replicates data. Here we do not expect to find differences, yet we see distorted latent factors being driven by the total UMIs. PCA was applied to 5,000 random genes.

in many cases there may be excess biological variability present as well. For example, if we wish to capture variability due to clusters of different cell types in a dimension reduction, we may wish to exclude biological variability due to cell cycle. Biological variability not accounted for by the sampling distribution may be accommodated by using a Dirichlet-multinomial likelihood, which is overdispersed relative to the multinomial. In practice, both the multinomial and Dirichlet-multinomial are computationally intractable, and may be approximated by the Poisson and negative binomial likelihoods, respectively (detailed derivations provided in Methods). We found that numerical convergence of GLM-PCA with negative binomial likelihood was unstable, due to the difficulty of estimating the dispersion parameter, so we focused on the simpler Poisson likelihood in our assessments. Intuitively, using Poisson instead of negative binomial implies we assume the biological variability is captured by the factor model and the unwanted biological variability is small relative to the sampling variability. We ran Poisson GLM-PCA on the technical and biological replicates negative control datasets and found it removed the spurious correlation between the first dimension and the total UMIs and fraction of zeros (Figure 4).

## 2.7 Deviance residuals provide fast approximation to GLM-PCA

One disadvantage of GLM-PCA is it depends on an iterative algorithm to obtain estimates for the latent factors, and is at least ten times slower than PCA. We therefore propose a fast approximation to GLM-PCA. When using PCA a common first step is to center and scale the data for each gene as z-scores. This is equivalent to the following procedure. First, specify a null model of constant gene expression across cells, assuming a normal distribution. Next, find the MLEs of its parameters for each gene (the mean and variance). Finally, compute residuals of the model as the z-scores (derivation provided in Methods). The fact that scRNA-Seq data are skewed, discrete, and possessing many zeros suggests the normality assumption may be inappropriate. Further, using z-scores does not account for variability in total UMIs across cells. Instead, we propose to replace the normal null model with a multinomial null model as a better match to the data generating mechanism. The analogs to z-scores under this model are called deviance and Pearson residuals. Mathematical formulae are presented in Methods. Use of multinomial residuals enables a fast transformation similar to z-scores that avoids difficulties of normalization and log-transformation by directly modeling counts. Additionally, this framework allows straightforward adjustment for covariates such as cell cycle signatures or batch labels.

## 2.8 Feature selection using deviance

Feature selection, or identification of informative genes, may be accomplished by ranking genes using the *deviance*, which quantifies how well each gene fits a null model of constant expression across cells. Unlike the competing highly variable or highly expressed genes methods, which are sensitive to normalization,

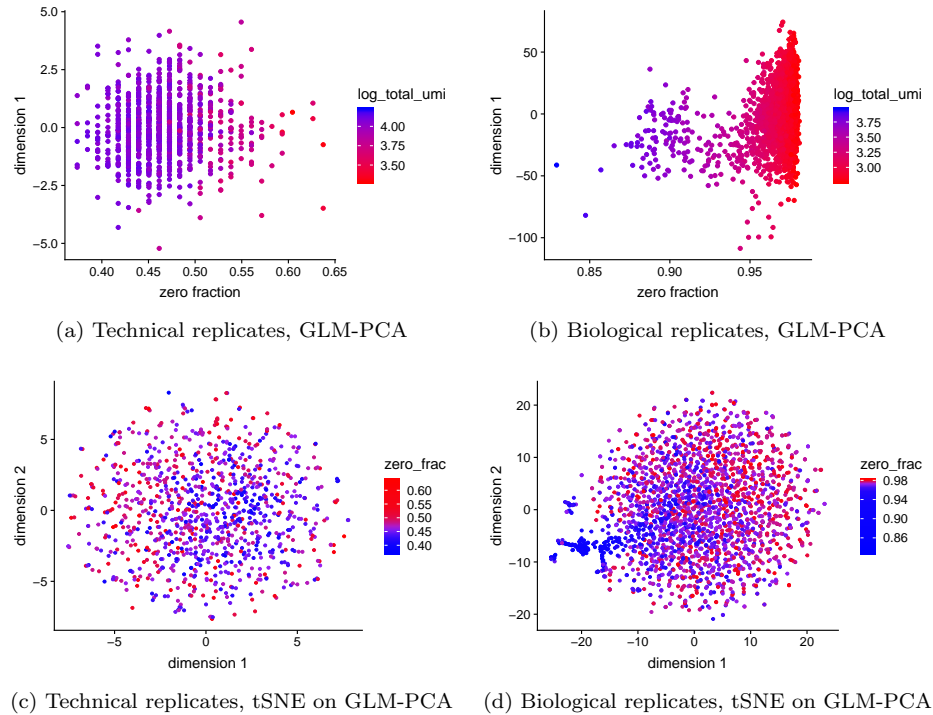


Figure 4: GLM-PCA dimension reduction is not affected by unwanted fraction of zeros variability and avoids false positive results. a) First GLM-PCA dimension (analogous to first principal component) plotted against the fraction of zeros for the technical replicates with colors representing the total UMIs. b) as a) but using biological replicates. c) Using the technical replicates, we applied t-distributed stochastic neighbor embedding (tSNE) with perplexity 30 to the top 50 GLM-PCA dimensions. The first two tSNE dimensions are shown with a blue to red color scale representing the fraction of zeros. d) as c) but for the biological replicates data. GLM-PCA using the Poisson approximation to the multinomial was applied to the same 5,000 random genes as in Figure 3.

ranking genes by deviance operates on raw UMI counts. An approximate multinomial deviance statistic can be computed in closed-form (formula provided in the Methods). We compared gene ranks for all three feature selection methods (deviance, highly expressed, and highly variable genes) on the 8eq dataset, which contained eight different known cell types. We found strong concordance between highly deviant genes and highly expressed genes (Spearman's rank correlation  $r = 0.9987$ ), while highly variable genes correlated weakly with both high expression ( $r = 0.3835$ ) and deviance ( $r = 0.3738$ ); see also Figure S2.

## 2.9 Multinomial models improve unsupervised clustering

Dimension reduction with GLM-PCA or its fast multinomial residuals approximation improved clustering performance over competing methods (Figure 5a). Feature selection by multinomial deviance was superior to highly variable genes (Figure 5b).

Using the two ground-truth datasets described in Section 2.1, we systematically compared the clustering performance of all combinations of previously described methods for normalization, feature selection, and dimension reduction. In addition, we compared against ZINB-WAVE since it also avoids requiring the user to pre-process and normalize the UMI count data (e.g. log transformation of CPM) and accounts for varying total UMIs across cells [27]. After obtaining latent factors, we used Seurat and k-means to infer clusters, and compared these to the known cell identities using Adjusted Rand Index (ARI, [37]). We varied the number of latent dimensions and number of clusters to assess robustness. Where possible, we used the same combinations of hyperparameters as [14] to facilitate comparisons to their extensive benchmarking (details are provided in Methods Section 4.6).

We compared the Seurat clustering performance of GLM-PCA (with Poisson approximation to multinomial) to running PCA on deviance residuals, which adhere more closely to the normal distribution than log-CPM. We found both of these approximate multinomial methods gave similar results on the 4eq dataset, and outperformed PCA on log-CPM z-scores. However, GLM-PCA outperformed the residuals method on the 8eq dataset. Also, performance on ZINB-WAVE factors degraded when the number of latent dimensions increased from 10 to 30, whereas GLM-PCA and its fast approximation with deviance residuals was robust to this change (Figure 5a). The performance of Pearson residuals was similar to that of deviance residuals (Figure S3).

Focusing on feature selection methods, deviance outperformed highly variable genes across both datasets and across dimension reduction methods (Figure 5b). Filtering by highly expressed genes led to similar clustering performance as deviance (Figure S3), because both criteria identified strongly overlapping gene lists for these data (Figure S2). The combination of feature selection with deviance and dimension reduction with GLM-PCA also improved clustering performance when k-means was used in place of Seurat (Figure S4). A complete table of results is publicly available (Section 5).

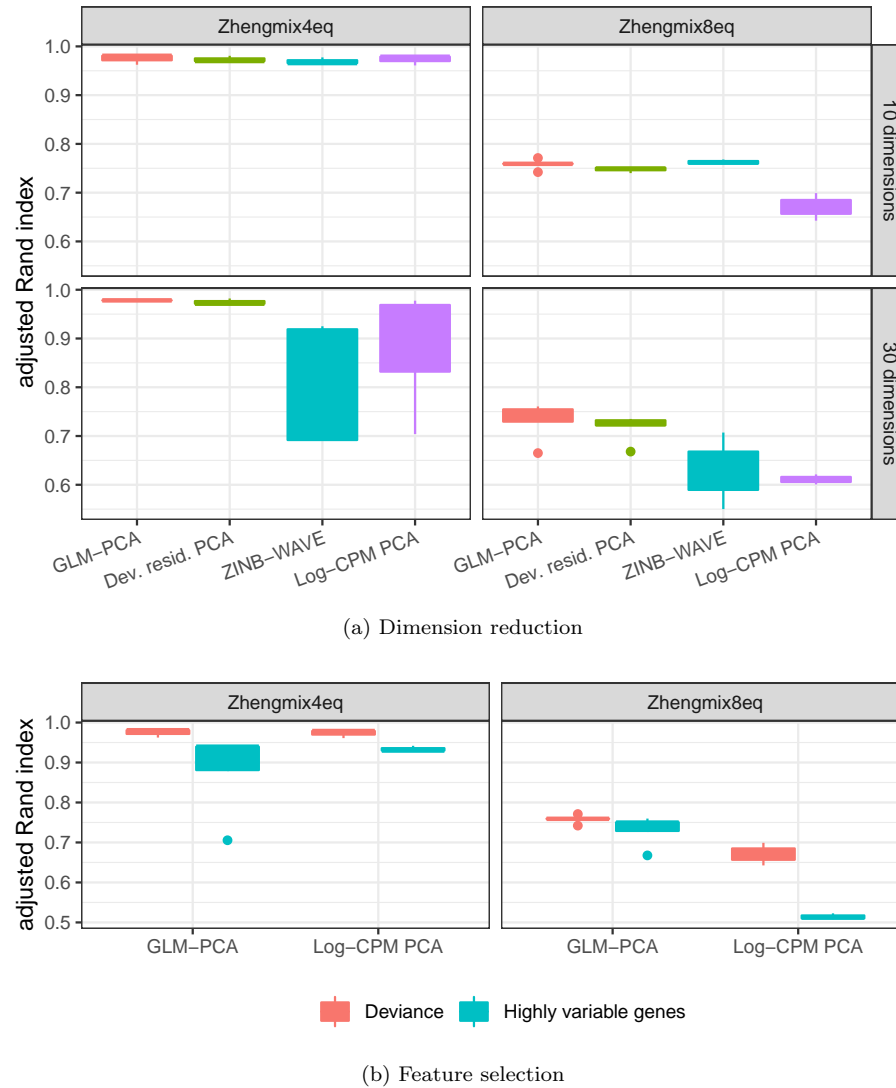


Figure 5: Dimension reduction with GLM-PCA and feature selection using deviance improves Seurat clustering performance. Each column represents a different ground-truth dataset from [14]. a) Comparison of dimension reduction methods based on the top 1,500 informative genes identified by approximate multinomial deviance. The Poisson approximation to the multinomial was used for GLM-PCA. Dev. resid. PCA: PCA on approximate multinomial deviance residuals. b) Comparison of feature selection methods. The top 1,500 genes identified by deviance and highly variable genes were passed to two different dimension reduction methods: GLM-PCA and PCA on log transformed CPM. Only results with the number of clusters within 25% of the true number are presented.

## 2.10 Computational efficiency of multinomial models

We measured time to convergence for reduction to two latent dimensions of GLM-PCA, ZINB-WAVE, PCA on log-CPM, PCA on deviance residuals, and PCA on Pearson residuals. Using the top 600 highly deviant genes, we subsampled the PBMC 68K dataset to 680, 6,800, and 68,000 cells. All methods scaled approximately linearly with increasing numbers of cells, but GLM-PCA was 23-63 times faster than ZINB-WAVE across sample sizes (Figure S5). Specifically, GLM-PCA processed 68,000 cells in less than seven minutes. The deviance and Pearson residuals methods exhibited speeds comparable to PCA: 9-26 times faster than GLM-PCA. We also timed dimension reduction of the 8eq dataset (3,994 cells) from 1,500 highly deviant genes to ten latent dimensions. PCA (with either log-CPM, deviance, or Pearson residuals) took 7 sec, GLM-PCA took 4.7 min, and ZINB-WAVE took 86.6 min.

## 3 Conclusions

We have outlined a statistical framework for analysis of scRNA-Seq data with UMI counts based on a multinomial model, providing effective and simple to compute methods for feature selection and dimension reduction. We found that UMI count distributions differ dramatically from read counts, are well-described by a multinomial distribution and are not zero-inflated. Log transformation of normalized UMI counts is detrimental, because it artificially exaggerates differences between zeros and all other values. For feature selection, or identification of informative genes, deviance is a more effective criterion than highly variable genes. Dimension reduction via GLM-PCA, or its fast approximation using residuals from a multinomial model, leads to better clustering performance than PCA on z-scores of log-CPM.

Although our methods were inspired by scRNA-Seq UMI counts, they may be useful for a wider array of data sources. Any high dimensional, sparse dataset where samples contain only relative information in the form of counts may conceivably be modeled by the multinomial distribution. Under such scenarios our methods are likely to be more effective than applying log-transformations and standard PCA. A possible example is microbiome data.

We have not addressed major topics in the scRNA-Seq literature such as pseudotime inference [38], differential expression [39], and spatial analysis [40]. However, the statistical ideas outlined here can also be used to improve methods in these more specialized types of analyses. In addition, adapting the GLM-PCA model to incorporate covariates such as batch labels or cell cycle signatures would be straightforward.

Our results have focused on (generalized) linear models for simplicity of exposition. Recently, several promising nonlinear dimension reductions for scRNA-Seq have been proposed. The variational autoencoder (VAE, a type of neural network) method scVI [41] utilizes a negative binomial likelihood in the decoder, while the encoder relies on log-normalized input data for numerical stability.

The Gaussian process method tGPLVM [42] models log-transformed counts. In both cases, we suggest replacing log-transformed values with deviance residuals to improve performance. Nonlinear dimension reduction methods may also depend on feature selection to reduce memory consumption and speed computation; here, our deviance method may be utilized as an alternative to high variability for screening informative genes.

The statistical approaches described here have not been validated against scRNA-Seq data without UMIs, such as SMART-Seq2 and other plate protocols [8], since non-UMI data contain PCR duplicates. To apply the ideas to these data, one would need to be able to infer the UMI counts for data with PCR replicates [7].

## 4 Methods

### 4.1 Multinomial Model for scRNA-Seq

Let  $y_{ij}$  be the observed UMI counts for cell or droplet  $i$  and gene or spike-in  $j$ . Let  $n_i = \sum_j y_{ij}$  be the total UMIs in the sample, and  $\pi_{ij}$  be the unknown true relative abundance of gene  $j$  in cell  $i$ . The random vector  $\mathbf{y}_i = (y_{i1}, \dots, y_{iJ})^\top$  with constraint  $\sum_j y_{ij} = n_i$  follows a multinomial distribution with density function

$$f(\mathbf{y}_i) = \binom{n_i}{y_{i1}, \dots, y_{iJ}} \prod_j \pi_{ij}^{y_{ij}}$$

Focusing on a single gene  $j$  at a time, the marginal distribution of  $y_{ij}$  is binomial with parameters  $n_i$  and  $\pi_{ij}$ . The marginal mean is  $E[y_{ij}] = n_i \pi_{ij} = \mu_{ij}$ , the marginal variance is  $\text{var}[y_{ij}] = n_i \pi_{ij} (1 - \pi_{ij}) = \mu_{ij} - \frac{1}{n_i} \mu_{ij}^2$ , and the marginal probability of a zero count is  $(1 - \pi_{ij})^{n_i} = \left(1 - \frac{\mu_{ij}}{n_i}\right)^{n_i}$ . The correlation between two genes  $j, k$  is

$$\text{cor}[y_{ij}, y_{ik}] = \frac{\sqrt{\pi_{ij} \pi_{ik}}}{\sqrt{(1 - \pi_{ij})(1 - \pi_{ik})}}$$

The correlation is induced by the sum to  $n_i$  constraint. As an extreme example, if there are only two genes ( $J = 2$ ), increasing the count of the first gene automatically reduces the count of the second gene since they must add up to  $n_i$  under multinomial sampling. This means when  $J = 2$  there is perfect anti-correlation between the gene counts which has nothing to do with biology. More generally, when either  $J$  or  $n_i$  is small, gene counts will be negatively correlated independent of biological gene-gene correlations, and it is not possible to analyze the data on a gene-by-gene basis (for example, by ranking and filtering genes for feature selection). Rather, comparisons are only possible between pairwise ratios of gene expression values [43]. Yet this type of analysis is difficult to interpret and computationally expensive for large numbers of genes (i.e. in high dimensions). Fortunately, under certain assumptions, more tractable approximations may be substituted for the true multinomial distribution.

First, note that if correlation is ignored, the multinomial may be approximated by  $J$  independent binomial distributions. Intuitively, this approximation will be reasonable if all  $\pi_{ij}$  are very small, which is likely to be satisfied for scRNA-Seq if the number of genes  $J$  is large, and no single gene constitutes the majority of mRNAs in the cell. If  $n_i$  is large and  $\pi_{ij}$  small, each binomial distribution can be further approximated by a Poisson with mean  $n_i\pi_{ij}$ . Alternatively, the multinomial can be constructed by drawing  $J$  independent Poisson random variables and conditioning on their sum. If  $J$  and  $n_i$  are large, the difference between the conditional, multinomial distribution and the independent Poissons becomes negligible. Since in practice  $n_i$  is large, the Poisson approximation to the multinomial may be reasonable [44, 45, 46, 47].

The multinomial model does not account for biological variability. As a result, an overdispersed version of the multinomial model may be necessary. This can be accommodated with the Dirichlet-multinomial distribution. Let  $\mathbf{y}_i$  be distributed as a multinomial conditional on the relative abundance parameter vector  $\pi_i = (\pi_{i1}, \dots, \pi_{iJ})^\top$ . If  $\pi_i$  is itself a random variable with symmetric Dirichlet distribution having shape parameter  $\alpha$ , the marginal distribution of  $\mathbf{y}_i$  is Dirichlet-multinomial. This distribution can itself be approximated by independent negative binomials. First, note that a symmetric Dirichlet random vector can be constructed by drawing  $J$  independent gamma variates with shape parameter  $\alpha$  and dividing by their sum. Suppose (as above) we approximate the conditional multinomial distribution of  $\mathbf{y}_i$  such that  $y_{ij}$  follows an approximate Poisson distribution with mean  $n_i\pi_{ij}$ . Let  $\lambda_{ij}$  be a collection of non-negative random variables such that  $\pi_{ij} = \frac{\lambda_{ij}}{\sum_j \lambda_{ij}}$ . We require that  $\pi_i$  follow a symmetric Dirichlet, which is accomplished by having  $\lambda_{ij}$  follow independent Gamma distributions with shape  $\alpha$  and mean  $n_i/J$ . This implies  $\sum_j \lambda_{ij}$  follows a Gamma with shape  $J\alpha$  and mean  $n_i$ . As  $J \rightarrow \infty$  this distribution converges to a point mass at  $n_i$ , so for large  $J$  (satisfied by scRNA-Seq),  $\sum_j \lambda_{ij} \approx n_i$ . This implies that  $y_{ij}$  approximately follows a conditional Poisson distribution with mean  $\lambda_{ij}$ , where  $\lambda_{ij}$  is itself a Gamma random variable with mean  $n_i/J$  and shape  $\alpha$ . If we then integrate out  $\lambda_{ij}$  we obtain the marginal distribution of  $y_{ij}$  as negative binomial with shape  $\alpha$  and mean  $n_i/J$ . Hence a negative binomial model for count data may be regarded as an approximation to an overdispersed Dirichlet-multinomial model.

Parameter estimation with multinomial models (and their binomial or Poisson approximations) is straightforward. First, suppose we observe replicate samples  $\mathbf{y}_i$ ,  $i = 1, \dots, I$  from the same underlying population of molecules, where the relative abundance of gene  $j$  is  $\pi_j$ . This is a null model because it assumes each gene has a constant expected expression level and there is no biological variation across samples. Regardless of whether one assumes a multinomial, binomial, or Poisson model, the maximum likelihood estimator (MLE) of  $\pi_j$  is  $\hat{\pi}_j = \frac{\sum_i y_{ij}}{\sum_i n_i}$  where  $n_i$  is the total count of sample  $i$ . In the more realistic case that relative abundances  $\pi_{ij}$  of genes vary across samples, the MLE is  $\hat{\pi}_{ij} = \frac{y_{ij}}{n_i}$ .

An alternative to the MLE is the maximum a posteriori (MAP) estimator.



Suppose a symmetric Dirichlet prior with concentration parameter  $\alpha_i$  is combined with the multinomial likelihood for cell  $i$ . The MAP estimator for  $\pi_{ij}$  is given by

$$\tilde{\pi}_{ij} = \frac{\alpha_i + y_{ij}}{J\alpha_i + n_i} = w_i \frac{1}{J} + (1 - w_i)\hat{\pi}_{ij}$$

where  $w_i = J\alpha_i/(J\alpha_i + n_i)$ , showing that the MAP is a weighted average of the prior mean that all genes are equally expressed ( $1/J$ ) and the MLE ( $\hat{\pi}_{ij}$ ). Compared to the MLE, the MAP biases the estimate toward the prior where all genes have the same expression. Larger values of  $\alpha_i$  introduce more bias, while  $\alpha_i \rightarrow 0$  leads to the MLE. If  $\alpha_i > 0$ , the smallest possible value of  $\tilde{\pi}_{ij}$  is  $\alpha_i/(J\alpha_i + n_i)$  rather than zero for the MLE. When there are many zeros in the data, MAP can stabilize relative abundance estimates at the cost of introducing bias.

## 4.2 Mathematics of distortion from log-normalizing UMIs

Suppose the true counts in cell  $i$  are given by  $x_{ij}$  for genes  $j = 1, \dots, J$ . Some of these may be zero, if a gene is not turned on in the cell. Knowing  $x_{ij}$  is equivalent to knowing the total number of transcripts  $t_i = \sum_j x_{ij}$  and the relative proportions of each gene  $\pi_{ij}$ , since  $x_{ij} = t_i \pi_{ij}$ . The total number of UMI counts  $n_i = \sum_j y_{ij}$  does not estimate  $t_i$ . However, under multinomial sampling, the UMI relative abundances  $\hat{\pi}_{ij} = \frac{y_{ij}}{n_i}$  are MLEs for the true proportions  $\pi_{ij}$ . Note that it is possible that  $\hat{\pi}_{ij} = 0$  even though  $\pi_{ij} > 0$ , indicating a dropout. Because  $\sum_j \hat{\pi}_{ij} = 1$  regardless of  $n_i$ , the use of multinomial MLEs is equivalent to the widespread practice of normalizing each cell by the total counts. Furthermore, the use of size factors  $s_i = n_i/m$  leads to  $\hat{\pi}_{ij} * m$  (if  $m = 10^6$  this is CPM).

Traditional bulk RNA-Seq experiments measured gene expression in read counts of many cells per sample rather than UMI counts of single cells. Gene counts from bulk RNA-Seq could thus range over several orders of magnitude. To facilitate comparison of these large numbers many bulk RNA-Seq methods have relied on a logarithm transformation. This enables interpretation of differences in normalized counts as fold changes on a relative scale. Prior to the use of UMIs, scRNA-Seq experiments also produced read counts with wide ranging values, and a log transform was again employed. However, with single cell data, more than 90% of the genes might be observed as exact zeros, and  $\log(0) = -\infty$  which is not useful for data analysis. UMI data also contain large numbers of zeros, but do not contain very large counts since PCR duplicates have been removed. Nevertheless, log transformation has been commonly used with UMI data as well.

The current standard is to transform the UMI counts as  $\log_2(c + \hat{\pi}_{ij} * m)$  where  $c$  is a pseudocount to avoid taking the log of zero, and typically  $c = 1$ . As before,  $m$  is some constant such as  $10^6$  for CPM. Finally, the data are centered and scaled so that the mean of each gene across cells is zero, and the standard deviation is one. This standardization of the data causes any

subsequent computation of distances or dimension reduction to be invariant to constant additive or multiplicative scaling. For example, under Manhattan distance  $d(x+c, y+c) = |x+c - (y+c)| = |x-y| = d(x, y)$ . In particular, using size factors such as CPM instead of relative abundances leads to a rescaling of the pseudocount, and use of any pseudocount is equivalent to replacing the MLE with the MAP estimator. Let  $k = c/m$  and  $\alpha_i = kn_i$ . Then the weight term in the MAP formula becomes  $w_i = Jk/(1 + Jk) = w$  which is constant across all cells  $i$ . Furthermore  $Jk = w/(1 - w)$ , showing that

$$\begin{aligned} \log_2(c + \hat{\pi}_{ij} * m) &= \log_2(k + \hat{\pi}_{ij}) + \log_2(m) \\ &= \log_2\left(\frac{w}{1-w} \frac{1}{J} + \hat{\pi}_{ij}\right) + \log_2(m) \\ &= \log_2\left(w \frac{1}{J} + (1-w)\hat{\pi}_{ij}\right) - \log_2(1-w) + \log_2(m) \\ &= \log_2(\tilde{\pi}_{ij}) + C \end{aligned}$$

Where  $C$  is a global constant that does not vary across cells or genes. For illustration, if  $c = 1$  and  $m = 10^6$  this is equivalent to assuming a prior where all genes are equally expressed and for cell  $i$ , a weight of  $w = J/(10^6 + J)$  is given to the prior relative to the MLE. Since the number of genes  $J$  is on the order of  $10^4$ , we have  $w \approx .01$ . The prior sample size for cell  $i$  is  $J\alpha_i = 10^{-6}Jn_i \approx .01*n_i$  where  $n_i$  is the data sample size. The standard transformation is therefore equivalent to using a weak prior to obtain a MAP estimate of the relative abundances, then log-transforming before dimension reduction.

In most scRNA-Seq datasets, the total number of UMIs  $n_i$  for some cells may be significantly less than the constant  $m$ . For these cells, the size factors  $s_i = n_i/m$  are less than one. Therefore, after normalization (dividing by size factor), the counts are scaled up to match the target size of  $m$ . Due to the discreteness of counts, this introduces a bias after log transformation, if the pseudocount is small (or equivalently, if  $m$  is large). For example, let  $c = 1$  and  $m = 10^6$  (CPM). If  $n_i = 10^4$  for a particular cell, we have  $s_i = .01$ . A raw count of  $y_{ij} = 1$  for this cell is normalized to  $1/.01 = 100$  and transformed to  $\log_2(1 + 100) = 6.7$ . For this cell, on the log scale there cannot be any values between zero and 6.7 because fractional UMI counts cannot be observed, and  $\log_2(1 + 0) = 0$ . Small pseudocounts and small size factors combined with log transform arbitrarily exaggerate the difference between a zero count and a small nonzero count. As previously shown, this scenario is equivalent to using MAP estimation of  $\pi_{ij}$  with a weak prior. To combat this distortion, one may attempt to strengthen the prior to regularize  $\hat{\pi}_{ij}$  estimation at the cost of additional bias, as advocated by [20]. An extreme case occurs when  $c = 1$  and  $m = 1$ . Here, the prior sample size is  $Jn_i$  so almost all the weight is on the prior. The transform is then  $\log_2(1 + \hat{\pi}_{ij})$ . But this function is approximately linear on the domain  $0 \leq \hat{\pi}_{ij} \leq 1$ . After centering and scaling, a linear transformation is vacuous.

To summarize, log transformation with a weak prior (small size factor, such as CPM) introduces strong artificial distortion between zeros and nonzeros,

while log transformation with a strong prior (large size factor) is roughly equivalent to not log transforming the data.

### 4.3 Generalized PCA

PCA minimizes the mean squared error (MSE) between the data and a low-rank representation, or embedding. Let  $y_{ij}$  be the raw counts and  $z_{ij}$  be the normalized and transformed version of  $y_{ij}$  such as centered and scaled log-CPM (z-scores). The PCA objective function is:

$$\min_{\mathbf{u}, \mathbf{v}} \sum_{i,j} (z_{ij} - \mathbf{u}'_i \mathbf{v}_j)^2$$

where  $\mathbf{u}_i, \mathbf{v}_j \in \mathbb{R}^L$  for  $i = 1, \dots, I$ ,  $j = 1, \dots, J$ . The  $\mathbf{u}_i$  are called factors or principal components and the  $\mathbf{v}_j$  are called loadings. The number of latent dimensions  $L$  controls the complexity of the model. Minimization of the MSE is equivalent to minimizing the Euclidean distance metric between the embedding and the data. It is also equivalent to maximizing the likelihood of a Gaussian model

$$z_{ij} \sim \mathcal{N}(\mathbf{u}'_i \mathbf{v}_j, \sigma^2)$$

If we replace the Gaussian model with a Poisson, which approximates the multinomial, we can directly model the UMI counts as

$$y_{ij} \sim \text{Poi}(n_i \exp\{\mathbf{u}'_i \mathbf{v}_j\})$$

or alternatively, in the case of overdispersion, we may approximate the Dirichlet-multinomial using a negative binomial likelihood

$$y_{ij} \sim \text{NB}(n_i \exp\{\mathbf{u}'_i \mathbf{v}_j\}; \phi_j)$$

We define the *linear predictor* as  $\eta_{ij} = \log n_i + \mathbf{u}'_i \mathbf{v}_j$ . It is clear that the mean  $\mu_{ij} = e^{\eta_{ij}}$  appears in both the Poisson and Negative Binomial model statements, showing that the latent factors interact with the data only through the mean. We may then estimate  $\mathbf{u}_i, \mathbf{v}_j$  (and  $\phi_j$ ) by maximizing the likelihood (in practice, adding a small L2 penalty to large parameter values improves numerical stability). A link function must be used since  $\mathbf{u}_i, \mathbf{v}_j$  are real valued whereas the mean of a Poisson or negative binomial must be positive. The total UMIs  $n_i$  term is used as an offset since no normalization has taken place; alternative size factors  $s_i$  such as those from *scran* [19] could be used in place of  $n_i$ . If the first element of each  $\mathbf{u}_i$  is constrained to equal 1, this induces a gene-specific intercept term in the first position of each  $\mathbf{v}_j$ , which is analogous to centering. Otherwise, the model is very similar to that of PCA; it is simply optimizing a different objective function. Unfortunately, MLEs for  $\mathbf{u}_i, \mathbf{v}_j$  cannot be expressed in closed form, so an iterative Fisher Scoring procedure is necessary. We refer to this model as GLM-PCA. Just as PCA minimizes MSE, GLM-PCA minimizes a generalization of MSE called the *deviance* [48]. While generalized PCA has been discovered before by [29], our implementation is novel in that

it allows for intercept terms, offsets, and non-canonical link functions. We also use a blockwise update for optimization which we found to be more numerically stable than that of [29]; we iterate over latent dimensions  $l$  rather than rows or columns. This technique is inspired by non-negative matrix factorization algorithms such as hierarchical alternating least squares and rank-one residue iteration, see [49] for a review.

As an illustration, consider GLM-PCA with the Poisson approximation to a multinomial likelihood. The objective function to be minimized is simply the overall deviance:

$$D = \sum_{i,j} y_{ij} \log \left( \frac{y_{ij}}{\mu_{ij}} \right) - (y_{ij} - \mu_{ij})$$

$$\log \mu_{ij} = \eta_{ij} = \log s_i + \mathbf{u}_i' \mathbf{v}_j = \log s_i + v_{j1} + \sum_{l=2}^L u_{il} v_{jl}$$

where  $s_i$  is a fixed size factor such as the total number of UMIs ( $n_i$ ). The optimization proceeds by taking derivatives with respect to the unknown parameters:  $v_{j1}$  is a gene-specific intercept term, and the remaining  $u_{il}, v_{jl}$  are the latent factors.

The GLM-PCA method is most concordant to the data generating mechanism since all aspects of the pipeline are integrated into a coherent model rather than being dealt with through sequential normalizations and transformations. The interpretation of the  $\mathbf{u}_i$  and  $\mathbf{v}_j$  vectors is the same as in PCA. For example, suppose we set the number of latent dimensions to two (i.e.  $L = 3$  to account for the intercept). We can plot  $u_{i2}$  on the horizontal axis and  $u_{i3}$  on the vertical axis for each cell  $i$  to visualize relationships between cells such as gradients or clusters. In this way, the  $\mathbf{u}_i$  and  $\mathbf{v}_j$  capture biological variability such as differentially expressed genes.

#### 4.4 Residuals and z-scores

Just as mean squared error can be computed by taking the sum of squared residuals under a Gaussian likelihood, the deviance is equal to the sum of squared *deviance residuals* [48]. Since deviance residuals are not well-defined for the multinomial distribution, we adopt the binomial approximation. The deviance residual for gene  $j$  in cell  $i$  is given by

$$r_{ij}^{(d)} = \text{sign}(y_{ij} - \hat{\mu}_{ij}) \sqrt{2y_{ij} \log \frac{y_{ij}}{\hat{\mu}_{ij}} + 2(n_i - y_{ij}) \log \frac{n_i - y_{ij}}{n_i - \hat{\mu}_{ij}}}$$

where under the null model of constant gene expression across cells,  $\hat{\mu}_{ij} = n_i \hat{\pi}_j$ . The deviance residuals are the result of regressing away this null model. An alternative to deviance residuals is the Pearson residual, which is simply the difference in observed and expected values scaled by an estimate of the standard

deviation. For the binomial, this is

$$r_{ij}^{(p)} = \frac{y_{ij} - \hat{\mu}_{ij}}{\sqrt{\hat{\mu}_{ij} - \frac{1}{n_i} \hat{\mu}_{ij}^2}}$$

According to the theory of generalized linear models (GLM), both types of residuals follow approximately a normal distribution with mean zero if the null model is correct [48]. Deviance residuals tend to be more symmetric than Pearson residuals. In practice, the residuals may not have mean exactly equal to zero, and may be standardized by scaling their gene-specific standard deviation just as in the Gaussian case.

The z-score is simply the Pearson residual where we replace the multinomial likelihood with a Gaussian (normal) likelihood, and use normalized values instead of raw UMI counts. Let  $q_{ij}$  be the normalized (possibly log-transformed) expression of gene  $j$  in cell  $i$  without centering and scaling. The null model is that the expression of the gene is constant across all cells:

$$q_{ij} \sim \mathcal{N}(\mu_j, \sigma_j^2)$$

The MLEs are  $\hat{\mu}_j = \frac{1}{I} \sum_i q_{ij}$ ,  $\hat{\sigma}_j^2 = \frac{1}{I} \sum_i (q_{ij} - \hat{\mu}_j)^2$ , and the z-scores equal the Pearson residuals  $z_{ij} = (q_{ij} - \hat{\mu}_j) / \hat{\sigma}_j$ .

## 4.5 Feature selection using deviance

Genes with constant expression across cells are not informative. Such genes may be described by the multinomial null model where  $\pi_{ij} = \pi_j$ . Goodness of fit to a multinomial distribution can be quantified using deviance, which is twice the difference in log-likelihoods comparing a saturated model to a fitted model. The multinomial deviance is a joint deviance across all genes and for this reason is not helpful for screening informative genes. Instead, one may use the binomial deviance as an approximation:

$$D_j = 2 \sum_i \left[ y_{ij} \log \frac{y_{ij}}{n_i \hat{\pi}_j} + (n_i - y_{ij}) \log \frac{(n_i - y_{ij})}{n_i (1 - \hat{\pi}_j)} \right]$$

A large deviance value indicates the model in question provides a poor fit. Those genes with biological variation across cells will be poorly fit by the null model and will have the largest deviances. By ranking genes according to their deviances, one may thus obtain highly deviant genes as an alternative to highly variable or highly expressed genes.

## 4.6 Systematic Comparison of Methods

We considered combinations of the following methods and parameter settings, following [14]. *Italics indicate methods proposed in this manuscript.* Feature selection: highly expressed genes, highly variable genes, and *highly deviant genes*. We did not compare against highly dropout genes because [14] found

this method to have poor downstream clustering performance for UMI counts and it is not as widely used in the literature. Number of genes: 60, 300, 1,500. Normalization, transformation, and dimension reduction: PCA on log-CPM z-scores, ZINB-WAVE [27], *PCA on deviance residuals*, *PCA on Pearson residuals*, and *GLM-PCA*. Number of latent dimensions: 10, 30. Clustering algorithm: k-means [50], Seurat [16]. Number of clusters: all values from 2-10, inclusive. Seurat resolution: 0.05, 0.1, 0.2, 0.5, 0.8, 1, 1.2, 1.5, 2.

## 5 Declarations

### Ethics approval and consent to participate

Not applicable.

### Consent for publication

Not applicable.

### Availability of data and material

All methods and assessments described in this manuscript are publicly available at <https://github.com/willtownes/scrna2019>.

### Competing interests

None declared.

### Funding

FWT is supported by NIH grant T32CA009337, SCH is supported by NIH grant R00HG009007, MJA is supported by an MGH Pathology Department startup fund, and RAI is supported by NIH grants R01HG005220, R01GM083084, and P41HG004059.

### Authors' contributions

SCH, MJA, and RAI identified the problem. FWT proposed, derived, and implemented the GLM-PCA model, its fast approximation using residuals, and feature selection using deviance. SCH, MJA and RAI provided guidance on refining the methods and evaluation strategies. FWT and RAI wrote the draft manuscript and revisions were suggested by SCH and MJA. All authors approved the final manuscript.

## Acknowledgements

The authors thank Keegan Korthauer, Jeff Miller, Linglin Huang, Alejandro Reyes, and Yered Pita-Juarez for valuable suggestions.

## References

- [1] Kalisky T, Oriel S, Bar-Lev TH, Ben-Haim N, Trink A, Wineberg Y, et al. A Brief Review of Single-Cell Transcriptomic Technologies. *Briefings in Functional Genomics*. 2018 Jan;17(1):64–76.
- [2] Svensson V, Vento-Tormo R, Teichmann SA. Exponential Scaling of Single-Cell RNA-Seq in the Past Decade. *Nature Protocols*. 2018 Apr;13(4):599–604.
- [3] Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, Goldman M, et al. Highly Parallel Genome-Wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell*. 2015 May;161(5):1202–1214.
- [4] Klein AM, Mazutis L, Akartuna I, Tallapragada N, Veres A, Li V, et al. Droplet Barcoding for Single-Cell Transcriptomics Applied to Embryonic Stem Cells. *Cell*. 2015 May;161(5):1187–1201.
- [5] Zheng GXY, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, et al. Massively Parallel Digital Transcriptional Profiling of Single Cells. *Nature Communications*. 2017 Jan;8:ncomms14049.
- [6] Dal Molin A, Di Camillo B. How to Design a Single-Cell RNA-Sequencing Experiment: Pitfalls, Challenges and Perspectives. *Briefings in Bioinformatics*. 2018 Jan;.
- [7] Qiu X, Hill A, Packer J, Lin D, Ma YA, Trapnell C. Single-Cell mRNA Quantification and Differential Analysis with Census. *Nature Methods*. 2017 Jan;advance online publication.
- [8] Picelli S, Björklund AK, Faridani OR, Sagasser S, Winberg G, Sandberg R. Smart-Seq2 for Sensitive Full-Length Transcriptome Profiling in Single Cells. *Nature Methods*. 2013 Nov;10(11):1096–1098.
- [9] Kolodziejczyk AA, Kim JK, Svensson V, Marioni JC, Teichmann SA. The Technology and Biology of Single-Cell RNA Sequencing. *Molecular Cell*. 2015 May;58(4):610–620.
- [10] Islam S, Zeisel A, Joost S, La Manno G, Zajac P, Kasper M, et al. Quantitative Single-Cell RNA-Seq with Unique Molecular Identifiers. *Nature Methods*. 2014 Feb;11(2):163–166.
- [11] Lun ATL, McCarthy DJ, Marioni JC. A Step-by-Step Workflow for Low-Level Analysis of Single-Cell RNA-Seq Data with Bioconductor. *F1000Research*. 2016 Oct;5:2122.

- [12] McCarthy DJ, Campbell KR, Lun ATL, Wills QF. Scater: Pre-Processing, Quality Control, Normalization and Visualization of Single-Cell RNA-Seq Data in R. *Bioinformatics*. 2017 Apr;33(8):1179–1186.
- [13] Andrews TS, Hemberg M. Identifying Cell Populations with scRNASeq. *Molecular Aspects of Medicine*. 2017 Jul;.
- [14] Duò A, Robinson MD, Sonesson C. A Systematic Performance Evaluation of Clustering Methods for Single-Cell RNA-Seq Data. *F1000Research*. 2018 Jul;7:1141.
- [15] Brennecke P, Anders S, Kim JK, Kołodziejczyk AA, Zhang X, Proserpio V, et al. Accounting for Technical Noise in Single-Cell RNA-Seq Experiments. *Nature Methods*. 2013 Nov;10(11):1093–1095.
- [16] Butler A, Hoffman P, Smibert P, Papalexi E, Satija R. Integrating Single-Cell Transcriptomic Data across Different Conditions, Technologies, and Species. *Nature Biotechnology*. 2018 Apr;.
- [17] Andrews TS, Hemberg M. Dropout-Based Feature Selection for scRNASeq. *bioRxiv*. 2018 May;p. 065094.
- [18] Hotelling H. Analysis of a Complex of Statistical Variables into Principal Components. *Journal of Educational Psychology*. 1933;24(6):417–441.
- [19] Lun AT, Bach K, Marioni JC. Pooling across Cells to Normalize Single-Cell RNA Sequencing Data with Many Zero Counts. *Genome Biology*. 2016;17:75.
- [20] Lun A. Overcoming Systematic Errors Caused by Log-Transformation of Normalized Single-Cell RNA Sequencing Data. *bioRxiv*. 2018 Aug;p. 404962.
- [21] Warton DI. Why You Cannot Transform Your Way out of Trouble for Small Counts. *Biometrics*. 2018 Mar;74(1):362–368.
- [22] Vallejos CA, Risso D, Scialdone A, Dudoit S, Marioni JC. Normalizing Single-Cell RNA Sequencing Data: Challenges and Opportunities. *Nature methods*. 2017 Jun;14(6):565–571.
- [23] Finak G, McDavid A, Yajima M, Deng J, Gersuk V, Shalek AK, et al. MAST: A Flexible Statistical Framework for Assessing Transcriptional Changes and Characterizing Heterogeneity in Single-Cell RNA Sequencing Data. *Genome Biology*. 2015 Dec;16:278.
- [24] Pierson E, Yau C. ZIFA: Dimensionality Reduction for Zero-Inflated Single-Cell Gene Expression Analysis. *Genome Biology*. 2015;16:241.
- [25] Liu S, Trapnell C. Single-Cell Transcriptome Sequencing: Recent Advances and Remaining Challenges. *F1000Research*. 2016 Feb;5.



- [26] Lin P, Troup M, Ho JWK. CIDR: Ultrafast and Accurate Clustering through Imputation for Single-Cell RNA-Seq Data. *Genome Biology*. 2017 Mar;18:59.
- [27] Risso D, Perraudeau F, Gribkova S, Dudoit S, Vert JP. ZINB-WaVE: A General and Flexible Method for Signal Extraction from Single-Cell RNA-Seq Data. *bioRxiv*. 2017 Nov;p. 125112.
- [28] Hicks SC, Townes FW, Teng M, Irizarry RA. Missing Data and Technical Variability in Single-Cell RNA-Sequencing Experiments. *Biostatistics*. 2018;19(4).
- [29] Collins M, Dasgupta S, Schapire RE. A Generalization of Principal Components Analysis to the Exponential Family. In: Dietterich TG, Becker S, Ghahramani Z, editors. *Advances in Neural Information Processing Systems 14*. MIT Press; 2002. p. 617–624.
- [30] Tung PY, Blischak JD, Hsiao CJ, Knowles DA, Burnett JE, Pritchard JK, et al. Batch Effects and the Effective Design of Single-Cell Gene Expression Studies. *Scientific Reports*. 2017 Jan;7:srep39921.
- [31] Ellefson JW, Gollihar J, Shroff R, Shivram H, Iyer VR, Ellington AD. Synthetic Evolutionary Origin of a Proofreading Reverse Transcriptase. *Science*. 2016 Jun;352(6293):1590–1593.
- [32] Shapiro E, Biezuner T, Linnarsson S. Single-Cell Sequencing-Based Technologies Will Revolutionize Whole-Organism Science. *Nature Reviews Genetics*. 2013 Sep;14(9):618–630.
- [33] Silverman JD, Roche K, Mukherjee S, David LA. Naught All Zeros in Sequence Count Data Are the Same. *bioRxiv*. 2018 Nov;p. 477794.
- [34] Pachter L. Models for Transcript Quantification from RNA-Seq. *arXiv:11043889 [q-bio, stat]*. 2011 Apr;.
- [35] Wagner F, Yan Y, Yanai I. K-Nearest Neighbor Smoothing for High-Throughput Single-Cell RNA-Seq Data. *bioRxiv*. 2018 Jan;p. 217737.
- [36] Van den Berge K, Perraudeau F, Soneson C, Love MI, Risso D, Vert JP, et al. Observation Weights Unlock Bulk RNA-Seq Tools for Zero Inflation and Single-Cell Applications. *Genome Biology*. 2018 Feb;19:24.
- [37] Hubert L, Arabie P. Comparing Partitions. *Journal of Classification*. 1985 Dec;2(1):193–218.
- [38] Trapnell C, Cacchiarelli D, Grimsby J, Pokharel P, Li S, Morse M, et al. The Dynamics and Regulators of Cell Fate Decisions Are Revealed by Pseudotemporal Ordering of Single Cells. *Nature Biotechnology*. 2014 Apr;32(4):381–386.

- [39] Soneson C, Robinson MD. Bias, Robustness and Scalability in Single-Cell Differential Expression Analysis. *Nature Methods*. 2018 Apr;15(4):255–261.
- [40] Svensson V, Teichmann SA, Stegle O. SpatialDE: Identification of Spatially Variable Genes. *Nature Methods*. 2018 Mar;.
- [41] Lopez R, Regier J, Cole MB, Jordan M, Yosef N. Bayesian Inference for a Generative Model of Transcriptome Profiles from Single-Cell RNA Sequencing. *bioRxiv*. 2018 Mar;p. 292037.
- [42] Verma A, Engelhardt B. A Robust Nonlinear Low-Dimensional Manifold for Single Cell RNA-Seq Data. *bioRxiv*. 2018 Oct;p. 443044.
- [43] Egozcue JJ, Pawłowsky-Glahn V, Mateu-Figueras G, Barceló-Vidal C. Isometric Logratio Transformations for Compositional Data Analysis. *Mathematical Geology*. 2003 Apr;35(3):279–300.
- [44] McDonald DR. On the Poisson Approximation to the Multinomial Distribution. *The Canadian Journal of Statistics / La Revue Canadienne de Statistique*. 1980;8(1):115–118.
- [45] Baker SG. The Multinomial-Poisson Transformation. *Journal of the Royal Statistical Society Series D (The Statistician)*. 1994;43(4):495–504.
- [46] Gopalan P, Hofman JM, Blei DM. Scalable Recommendation with Poisson Factorization. *arXiv:13111704 [cs, stat]*. 2013 Nov;.
- [47] Taddy M. Distributed Multinomial Regression. *The Annals of Applied Statistics*. 2015 Sep;9(3):1394–1414.
- [48] Agresti A. *Foundations of Linear and Generalized Linear Models*. John Wiley & Sons; 2015.
- [49] Kim J, He Y, Park H. Algorithms for Nonnegative Matrix and Tensor Factorizations: A Unified View Based on Block Coordinate Descent Framework. *Journal of Global Optimization*. 2014 Feb;58(2):285–319.
- [50] Hartigan JA, Wong MA. Algorithm AS 136: A K-Means Clustering Algorithm. *Journal of the Royal Statistical Society Series C (Applied Statistics)*. 1979;28(1):100–108.

## Supplemental Figures

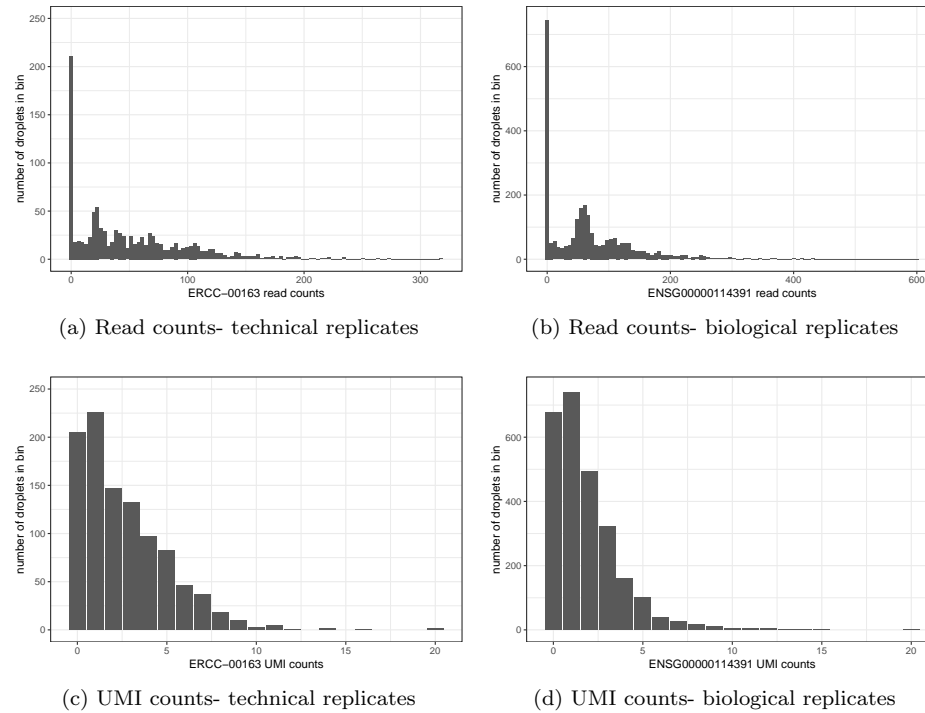


Figure S1: Comparing read counts and UMI counts sampling distribution from technical and biological replicates negative control datasets. a) Read count distribution for spike-in ERCC-00163 across technical replicates. b) Read count distribution for gene ENSG00000114391 across biological replicates (purified monocytes). c) as a) but without PCR duplicates. d) as b) but without PCR duplicates.

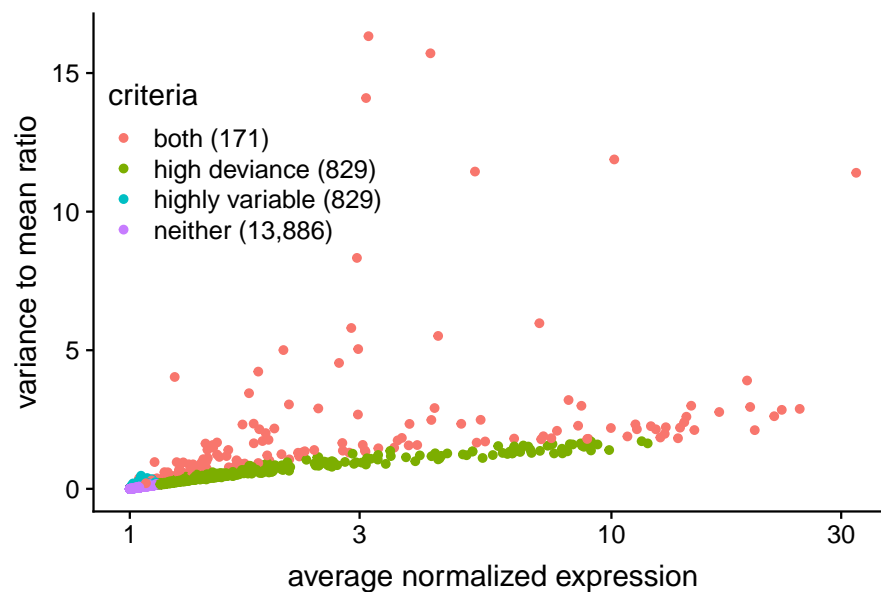


Figure S2: Comparison of top 1,000 genes selected as most informative in the Zheng 8eq dataset. The variance to mean ratio for each gene is plotted against the average expression. Counts were normalized using scran [19]. Colors represent genes that are in the top 1,000 ranked by variability (blue, red) and top 1,000 ranked by approximate multinomial deviance (green, red). Red indicates genes identified by both criteria, while purple indicates genes identified by neither criteria. Note that highly expressed genes have large values on the horizontal axis. The number of genes in each category is shown in parentheses.

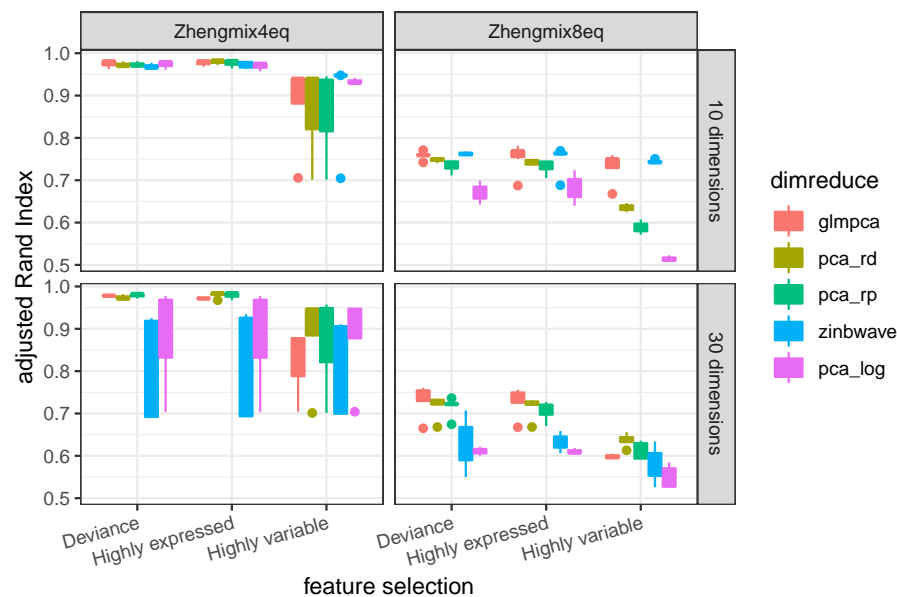


Figure S3: Comparison of Seurat clustering performance for all dimension reduction and feature selection methods on ground-truth datasets from [14]. The number of informative genes was fixed at 1,500. The Poisson approximation to the multinomial was used for GLM-PCA. Only results with the number of clusters within 25% of the true number are presented. Abbreviations: dimreduce: dimension reduction method, pca\_rd: PCA on deviance residuals, pca\_rp: PCA on Pearson residuals, pca\_log: PCA on log-CPM.

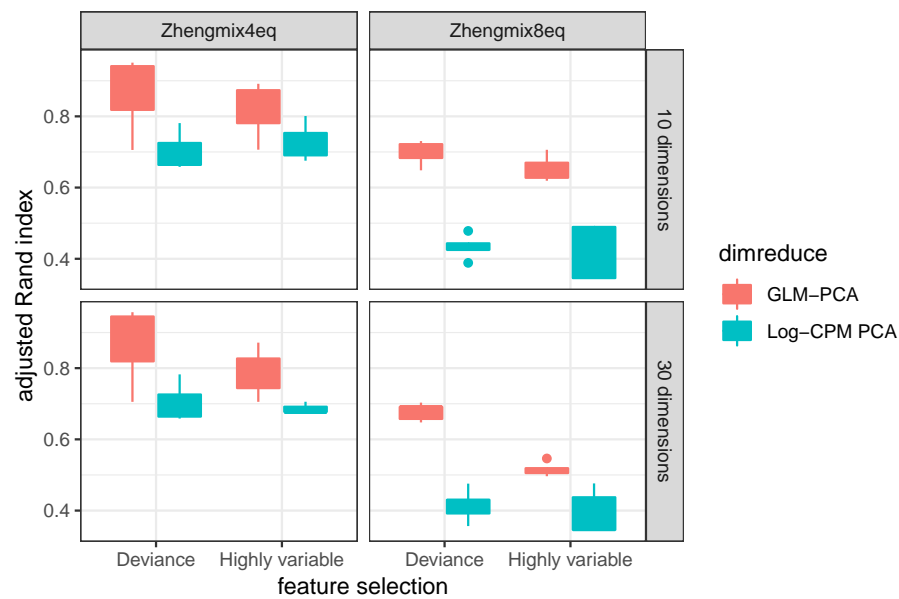


Figure S4: Dimension reduction with GLM-PCA and feature selection using deviance improves k-means clustering performance. Each column represents a different ground-truth dataset from [14]. The top 1,500 informative genes were identified by approximate multinomial deviance and highly variable genes. The Poisson approximation to the multinomial was used for GLM-PCA. Only results with the number of clusters within 25% of the true number are presented.

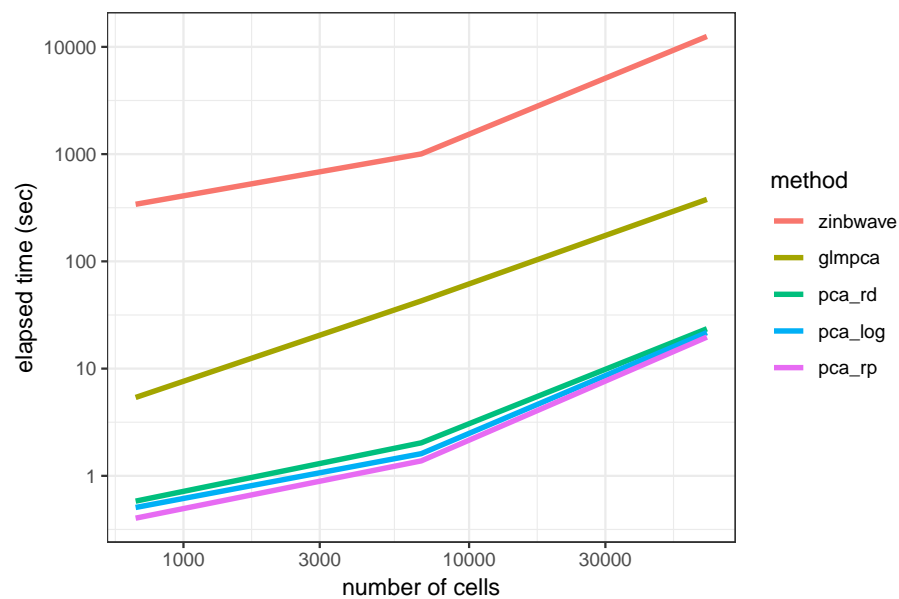


Figure S5: Computational speed comparison of dimension reduction methods GLM-PCA (glmpca), ZINB-WAVE (zinbwave), PCA on deviance residuals (pca\_rd), PCA on Pearson residuals (pca\_rp), and PCA on log-CPM (pca\_log).