

# A COMPOSITIONAL MODEL TO ASSESS EXPRESSION CHANGES FROM SINGLE-CELL RNA-SEQ DATA

XIUYU MA, AND CHRISTINA KENDZIORSKI, AND MICHAEL A. NEWTON

## 1. INTRODUCTION

The ability to measure genome-wide gene expression at single-cell resolution has accelerated the pace of biological discovery[1]. Overcoming data analysis challenges caused by the scale and unique variation properties of single-cell data will surely fuel further advances in immunology (cite), developmental biology (cite), and cancer (cite)(unsure about which paper to cite). Computational tools and statistical methodologies created for data of lower-resolution (e.g. bulk RNA-seq) or lower dimension (e.g. flow cytometry) guide our response to the data science demands of new measurement platforms, but they are not adequate for efficient knowledge discovery in this rapidly advancing domain[2].

An important feature of single-cell studies that could be leveraged better statistically is the fact that cells populate distinct, identifiable subtypes determined by lineage history, epigenetic state, the activity of various transcriptional programs (e.g. burst states), or other distinguishing factors. Lots of efforts have been made to clustering cells into different cell subtypes, SC3[3], CIDR[4] and ZIFA[5]. Whether or not a determination of cellular subtypes and their frequencies is a task of interest in a given application, we hypothesize that such subtype information may be injected into other inferences in order to improve their operating characteristics.

Assessing the magnitude and statistical significance of changes in gene expression associated with different cellular conditions has been a central statistical problem in genomics for which new tools specific to the single-cell RNAseq data structure have been deployed: MAST[6], DESEQ2[7], SCDD[8], etc. These tools respond to scRNAseq characteristics, such as high prevalence of zero counts and gene-level multimodality, but none takes explicit advantage of cellular subtype information. We present a simple procedure and supporting theoretical analyses for this purpose. A notable technical innovation is a new prior distribution over pairs of multinomial probability vectors that conveys both marginal Dirichlet conjugacy as well as dependence induced through sharp equalities on aggregated subtype probabilities, which turns out to be key in formulating the posterior probability of changes in expression distributions between conditions.

---

DEPARTMENT OF BIostatISTICS AND MEDICAL INFORMATICS, UW MADISON,  
TECHNICAL REPORT TR\*\*\*-V1, MAY \*\*, 2018.

In our compositional model, subtype that inferred from whole genome are fed into gene level expression. We utilize the mixture of subtypes to characterize transcripts profile and identify differential distributed genes across conditions in an scRNA-seq experiment. Simulation study suggests that the method provides improved power and precision for identifying differentially distributed genes. Performance on real data has been investigated through ten previously published experimental data from conquer[9]. We also obtained asymptotic properties of posterior inference.

## 2. MODELING

**2.1. Data structure, sampling model, and parameters.** In modeling scRNASeq data, we imagine that each cell  $c$  falls into one of  $K > 1$  classes, which we think of as subtypes or subpopulations of cells. For notation,  $z_c = k$  means that cell  $c$  happens to be of subtype  $k$ , with the vector  $z = (z_c)$  recording the state of all sampled cells. Knowledge of this class structure prior to measurement is not required, as it will be inferred as necessary from available genomic data. We also assume that cells arise from multiple experimental conditions, such as by treatment-control status or some other factor measured at the cell level, and we present our development for the special case of two conditions. Notationally,  $y = (y_c)$  records the experimental condition, say  $y_c = 1$  or  $y_c = 2$  initially: extensions to multiple conditions are discussed in section 6. Let's say condition  $j$  measures  $n_j = \sum_c 1[y_c = j]$  cells, and in total we have  $n = n_1 + n_2$  cells in the analysis. Further let  $t_k^j = \sum_c 1[y_c = j, z_c = k]$  denote the number of cells of subtype  $k$  in condition  $j$ ; we'll infer something about these counts using genome-wide data. As for molecular data, the normalized expression of gene  $g$  in cell  $c$ , say  $X_{g,c}$ , is one entry in a typically large GENES by CELLS data matrix  $X$ . Summing up, the data structure entails an expression matrix  $X$ , a treatment label vector  $y$ , and a vector  $z$  of latent subtype labels.

We treat subtype counts in the two conditions,  $t^1 = (t_1^1, t_2^1, \dots, t_K^1)$  and  $t^2 = (t_1^2, t_2^2, \dots, t_K^2)$ , as independent multinomial vectors, reflecting the common, two-condition experimental design. Explicitly,

$$t^1 \sim \text{Multinomial}_K(n_1, \phi) \quad \text{and} \quad t^2 \sim \text{Multinomial}_K(n_2, \psi)$$

for probability vectors  $\phi = (\phi_1, \phi_2, \dots, \phi_K)$  and  $\psi = (\psi_1, \psi_2, \dots, \psi_K)$  that characterize the populations of cells from which the  $n$  observed cells are sampled. Our working hypothesis, referred to as the *compositional model*, is that any differences in the distribution of expression  $X_{g,c}$  between  $y_c = 1$  and  $y_c = 2$  (i.e., any condition effects) are attributable to differences between the conditions in the underlying composition of cell types; i.e., owing to  $\phi \neq \psi$ . We reckon that cells of any given subtype  $k$  will present data according to a distribution reflecting technical and biological variation specific to that class of cells, regardless of the condition the cell finds itself in. Some care is needed in this, as an overly broad cell subtype (e.g. *epithelial cells*) could have further subtypes that show differential response to some treatment, for example, and so cellular condition (treatment) would then affect the distribution of expression data within the subtype, which is contrary to our working

hypothesis. On the other hand, we could then refine the subtype definition to allow more population classes  $K$  in order to mitigate that problem. There's a risk in this approach if  $K$  approaches  $n$  (i.e., every cell is its own type). In spite of this theoretical possibility, the data sets often encountered appear not to display this phenomenon, even when relatively flexible in the within-subtype expression distribution. We revisit the issue in discussion section, but for now proceed assuming that cellular condition affects the composition of subtypes but not the distribution of expression within a subtype.

With this compositional model, let  $f_{g,k}(x)$  denote the sampling distribution of expression measurement  $X_{g,c}$  assuming that cell  $c$  is from subtype  $k$ . Then in the two cellular conditions, the marginal distributions over subtypes are

$$f_g^1(x) = \sum_{k=1}^K \phi_k f_{g,k}(x) \quad \text{and} \quad f_g^2(x) = \sum_{k=1}^K \psi_k f_{g,k}(x).$$

We say that gene  $g$  is *differentially distributed*, denote  $\text{DD}_g$ , if  $f_g^1(x) \neq f_g^2(x)$  for some  $x$ , and otherwise it is equivalently distributed ( $\text{ED}_g$ ). Motivated by findings from bulk RNAseq data analysis, we further set each  $f_{g,k}$  to have a Negative Binomial form, say with mean  $\mu_{g,k}$  and shape parameter  $\alpha_g$  ([10]; add more citations here). This choice proves to be effective in our numerical experiments though it is not critical to the modeling formulation.

We seek a useful methodology to prioritize genes for evidence of  $\text{DD}_g$ . Interestingly, even if we have evidence for condition effects on the subtype frequencies, it does not follow that a given gene will have  $f_g^1 \neq f_g^2$ ; that depends on whether or not the subtypes show the right pattern of *differential expression* at  $g$ , to use the standard terminology from bulk RNAseq. For example, if two subtypes have different frequencies between the two conditions ( $\phi_1 \neq \psi_1$  and  $\phi_2 \neq \psi_2$ ) but the same aggregate frequency ( $\phi_1 + \phi_2 = \psi_1 + \psi_2$ ), and also if  $\mu_{g,1} = \mu_{g,2}$  then, other things being equal,  $f_g^1 = f_g^2$  even though  $\phi \neq \psi$ . Simply, a gene that does not distinguish two subtypes will also not distinguish the cellular conditions if those subtypes appear in the same aggregate frequency in the two conditions, regardless of changes in the individual subtype frequencies. We formalize this idea in order that our methodology has the necessary functionality. First, consider the parameter space

$$\Theta = \{(\phi, \psi, \mu, \sigma)\}$$

where  $\phi = (\phi_1, \phi_2, \dots, \phi_K)$  and  $\psi = (\psi_1, \psi_2, \dots, \psi_K)$ , as before, where  $\mu = \{\mu_{g,k}\}$ , all the subtype-and-gene-specific expected values, and where  $\sigma = \{\sigma_g\}$  holds all the gene-specific Negative binomial shape parameters. We define special subsets of  $\Theta$  using partitions of the  $K$  cell subtypes. A single partition, say  $\pi$ , is a set of mutually exclusive and exhaustive blocks,  $b$ , say, each a subset of  $\{1, 2, \dots, K\}$ , and we write  $\pi = \{b\}$ . We recall that the set  $\Pi$  containing all partitions  $\pi$  of  $\{1, 2, \dots, K\}$  has cardinality that grows rapidly with  $K$ . We'll carry along an example involving  $K = 7$  cell types, and one three-block partition taken from the set of 877 possible partitions of  $\{1, 2, \dots, 7\}$  (Figure 1).



FIGURE 1. Proportions of  $K = 7$  cellular subtypes in different conditions. Aggregated proportions of subtype 1 and 2, subtype 3, 4, 5, and subtype 6,7 remained same across conditions while proportion of individual subtype changed. A gene, for example, that shows no differential expression between subtypes 1 and 2, and also none among types 3,4, and 5, and none between 6 and 7 has the same marginal distribution between the two conditions.

For any partition  $\pi = \{b\}$  we have aggregate subtype frequencies

$$\Phi_b = \sum_{k \in b} \phi_k \quad \text{and} \quad \Psi_b = \sum_{k \in b} \psi_k.$$

We'll also use the notation  $\Phi_\pi = \{\Phi_b : b \in \pi\}$  and similarly for  $\Psi_\pi$ . As long as  $\pi$  is not the most refined partition, the mapping from  $(\phi, \psi)$  to  $(\Phi_\pi, \Psi_\pi)$  is many-to-one (Figure 1) Define

$$A_\pi = \{\theta \in \Theta : \Phi_b = \Psi_b \forall b \in \pi\}.$$

and

$$M_{g,\pi} = \{\theta \in \Theta : \mu_{g,k} = \mu_{g,k'} \iff k, k' \in b, b \in \pi\}.$$

Indeed, these are precisely the structures needed to address differential distribution  $DD_g$  (and its complement, equivalent distribution,  $ED_g$ ) at a given gene  $g$ :

**Theorem 1.** *Let  $C_{g,\pi} = A_\pi \cap M_{g,\pi}$ . For distinct partitions  $\pi_1, \pi_2$ ,  $C_{g,\pi_1} \cap C_{g,\pi_2} = \emptyset$ . Further, at any gene  $g$ , equivalent distribution is*

$$ED_g = \bigcup_{\pi \in \Pi} C_{g,\pi}.$$

This representation is central to our empirical Bayes posterior probabilities,  $P(DD_g|X, y)$ , which we develop below, and which we use to score differential distribution per gene. We

require a natural and simplifying prior assumption:  $A_\pi$  and  $M_{g,\pi}$  are independent. (I.e., proportions  $(\phi, \psi)$  are uninformative about the mean expression  $\{\mu_{g,i}\}$ ). Then,

$$(1) \quad 1 - P(\text{DD}_g|X, y) = P(\text{ED}_g|X, y) = \sum_{\pi \in \Pi} P(A_\pi|X, y) P(M_{g,\pi}|X).$$

The idea is by allowing genome-wide information to inform the posterior of  $A_\pi$ , then get added benefit from the gene-level data, which primarily drives the posterior of  $M_{g,\pi}$ . It's a very specific form of information sharing that leverages the compositional structure of single-cell data.

**2.2. Method structure.** Our approach take transcripts processed by normalizing methods (e.g. SCnorm [11]). The workflow contains two parts, clustering cells into subtypes and posterior inference on distributional change. Recall subtype is a group of cells with distributions of transcripts that are specific to this group, regardless which condition the cells is from. Thus identification process is blind to conditions and can be done by clustering upon similarities between cells.

---

**Algorithm 1**


---

**Input:** Expression data  $X$ ; condition labels  $y$ ; number of cell subtypes  $K$

**Output:** per gene posterior probability of differential distribution

---

```

1: procedure SCDDBOOST0( $X, y, K$ )
2:   function CLUSTERING( $X, K$ )
3:     distance matrix:  $\text{Dist}(X) \leftarrow$  pairwise distances between cells (columns of  $X$ )
4:     cell clustering:  $\hat{z} \leftarrow$  labeled partition of cells, computed from  $\text{Dist}(X)$  and  $K$ 
5:   function POST( $X, y, \hat{z}$ )
6:     subtype differential expression:  $\forall \pi, P(M_{g,\pi}|X, \hat{z}) \leftarrow$  using EBSeq[10]
7:     cell frequency changes:  $\forall \pi, P(A_\pi|y, \hat{z}) \leftarrow$  using Double Dirichlet prior
8:     posterior probability:  $P(\text{ED}_g|X, y) \leftarrow \sum_{\pi} P(M_{g,\pi}|X, \hat{z}) P(A_\pi|y, \hat{z})$ 
9:   return  $P(\text{DD}_g|X, y) = 1 - P(\text{ED}_g|X, y)$ 
```

---

One advantage of our approach is that the posterior inference can be incorporate with different clustering methods. With the development of technology, new clustering methods with higher accuracy may be discovered and we can substitute those methods with clustering function in our procedure. Relying on appropriate partition results, we gain power in identifying DD genes. No matter what clustering method we are using, the essential idea is that the mixture structure is estimated using the whole genome information rather than at each gene solely using information of that gene to estimate a gene specific mixture structure. Due to this reason, our model is more capable of capturing characteristic of scRNA seq data than scDD and we name our approach scDDboost.

There are many single cell clustering methods(e.g. SC3[3], CIDR[4] and ZIFA[5]) that could utilized in the workflow. However, for the integrity of our work and speeding in

computation, we present another efficient and effective clustering procedure. We pool cells from two conditions. At each gene level, we do a Poisson-Gamma model extended from modal clustering[12]. After gene level clustering, we use the cluster-based similarity partition algorithm (CSPA[13]). For each individual clustering result, a binary similarity matrix is constructed from the corresponding cell labels: if two cells belong to the same cluster, their similarity is 1; otherwise their similarity is 0. We obtain a consensus similarity matrix  $M_1$  by averaging all similarity matrices of individual clusterings. Another distance matrix  $M_2$  calculated by Pearson distance between cells. A final similarity matrix is obtained by weighted combining  $M_1$  and  $M_2$ . Cells are classified into subtypes by K-medoids based on the final similarity matrix.

After identification of subtypes, the following two steps use empirical Bayes to provide posterior inference on patterns of differential expression( $M_{g,\pi}$ ) and aggregated proportions of subtypes( $A_\pi$ ), where  $P(M_{g,\pi}|X)$  is done in EBSeq[10] and we present details of calculating  $P(A_\pi|X)$  in next section. Combining those components, final posterior probabilities are obtained.

### 3. EMPIRICAL BAYES

**3.1. Double Dirichlet prior.** Here we describe a prior  $p(\phi, \psi)$  that is conjugate to multinomial sampling but that also enables downstream gene-specific inferences about differential distribution when certain cell types do not differ in their expression distributions.

For our purposes, the prior will have a spike-slab structure that mixes over distinct patterns of equality of  $\pi$ -associated accumulated probabilities:

$$p(\phi, \psi) = \sum_{\pi \in \Pi} P(A_\pi) p(\phi, \psi | A_\pi)$$

Upon setting up a prior  $p(\phi, \psi)$  that can mix over structures  $A_\pi$ , we can obtain posterior inference  $P(A_\pi | t^1, t^2)$

Initially, the multitude of  $P(A_\pi)$ 's will be preset constants. To complete the prior specification  $p(\phi, \psi)$ , consider further scalars  $\alpha_k > 0$  for each class  $k$  and  $\beta_b > 0$  for each potential block  $b$ . (Extending the notational convention,  $\alpha_b$  is the vector of  $\alpha_k$  for  $k \in b$ , and  $\beta_\pi$  is the vector of  $\beta_b$  for  $b \in \pi$ .) For any block  $b$  consider conditional probabilities

$$\tilde{\phi}_b = \frac{\phi_b}{\Phi_b} \quad \tilde{\psi}_b = \frac{\psi_b}{\Psi_b}$$

which indicate the conditional probability of each class  $k$  given that the cell is of one of the types in  $b$ . Assume that conditional upon  $A_\pi$ ,

$$\Phi_\pi \sim \text{Dirichet}_{N(\pi)}[\beta_\pi]$$

where  $N(\pi)$  is the number of blocks  $b$  in  $\pi$ , and further that accumulated probabilities are the same between the two source conditions:  $\Phi_\pi = \Psi_\pi$ . Finally, assume that for each

$b \in \pi$ ,

$$\tilde{\phi}_b, \tilde{\psi}_b \sim_{\text{i.i.d.}} \text{Dirichlet}_{N(b)}[\alpha_b]$$

where  $N(b)$  is the number of cell types in block  $b$ . In other words, if  $A_\pi$  is the active structure, then accumulated probability vectors  $\Phi_\pi$  and  $\Psi_\pi$  are equal between the two source conditions, though the sub-block class-specific rates  $\phi_k$  and  $\psi_k$  may differ, as would (re-normalized) independent Dirichlet-distributed vectors. Taken together,

$$p(\phi, \psi | A_\pi) = p(\Phi_\pi, \Psi_\pi | A_\pi) \prod_{b \in \pi} [p(\tilde{\phi}_b) p(\tilde{\psi}_b)]$$

with

$$p(\Phi_\pi, \Psi_\pi | A_\pi) = \frac{\Gamma(\sum_{b \in \pi} \beta_b)}{\prod_{b \in \pi} \Gamma(\beta_b)} \left[ \prod_{b \in \pi} \Phi_b^{\beta_b - 1} \right] \mathbb{1}[\Phi_\pi = \Psi_\pi]$$

and

$$p(\tilde{\phi}_b) = \frac{\Gamma(\sum_{k \in b} \alpha_k)}{\prod_{k \in b} \Gamma(\alpha_k)} \prod_{k \in b} \tilde{\phi}_k^{\alpha_k - 1}, \quad p(\tilde{\psi}_b) = \frac{\Gamma(\sum_{k \in b} \alpha_k)}{\prod_{k \in b} \Gamma(\alpha_k)} \prod_{k \in b} \tilde{\psi}_k^{\alpha_k - 1}.$$

**3.2. Predictive probabilities:** For notation, we use  $\phi_b$  for the vector of values  $\phi_k$  for  $k \in b$ , and similarly for  $\psi_b$ . Analogously,  $\Phi_\pi$  and  $\Psi_\pi$  are vectors of accumulated class probabilities  $\phi_b$  and  $\psi_b$  for all  $b \in \pi$ , respectively.

In order to get the posterior probability  $p(A_\pi | t^1, t^2)$ , we need to calculate

$$\begin{aligned} p(A_\pi | t^1, t^2) &\propto p(A_\pi, t^1, t^2) = \int_{A_\pi} p(t^1, t^2 | \phi, \psi) p(\phi, \psi) d\phi d\psi \\ &= \sum_{\pi' \in \Pi} \int_{A_\pi} p(t^1, t^2 | \phi, \psi) p(\phi, \psi | A_{\pi'}) p(A_{\pi'}) d\phi d\psi \end{aligned}$$

For simplicity of notation, let  $w(\pi_1, \pi_2) = \int_{A_{\pi_1}} p(t^1, t^2 | \phi, \psi) p(\phi, \psi | A_{\pi_2}) p(A_{\pi_2}) d\phi d\psi$ , then  $p(A_\pi | t^1, t^2) \propto \sum_{\pi' \in \Pi} w(\pi, \pi')$ . To calculate component  $w(\pi, \pi')$ , recall refinement and coarseness relationship between partitions, we say a partition  $\tilde{\pi}$  is a refinement of another partition  $\pi$  if  $\forall b \in \pi$  there exists  $s \subset \tilde{\pi}$  such that  $\bigcup_{b' \in s} b' = b$ . We say  $\pi$  is a coarseness of  $\tilde{\pi}$  when  $\tilde{\pi}$  refines  $\pi$ . we have following theorem

**Theorem 2.** *If  $\pi'$  is a refinement of  $\pi$  then  $w(\pi, \pi') = w(\pi', \pi')$  otherwise  $w(\pi, \pi') = 0$*

Consequently, let  $RF(\pi)$  be the collection of finer partition of  $\pi$ , we have the posterior probability:

$$p(A_\pi | t^1, t^2) \propto \sum_{\pi' \in RF(\pi)} w(\pi', \pi')$$

Using the Dirichlet-Multinomial conjugacy and the collapsing property of these distributions ([14]), we get closed formulas for the predictive probability of cell-type counts  $t^1$  and  $t^2$ . Fixing  $\pi$ , let  $t_b^j = \sum_{k \in b} t_k^j$ , for cell conditions  $j = 1, 2$ , record the total numbers of cells accumulated over all types in block  $b$ . And following our notation convention,  $t_\pi^j$  is the vector of these counts over  $b \in \pi$ . From the prior and model structure

$$w(\pi, \pi) = p(t^1 | t_\pi^1) p(t^2 | t_\pi^2) p(t_\pi^1, t_\pi^2 | A_\pi) p(A_\pi).$$

Conditional independence of  $t^1$  and  $t^2$  given the block-level totals  $t_\pi^1$  and  $t_\pi^2$  on  $A_\pi$  reflects the possible differential class proportion structure within blocks but between cell conditions. For either cellular group  $j = 1, 2$ , we find, after some simplification, the following Dirichlet-Multinomial masses:

$$(2) \quad p(t^j | t_\pi^j) = \prod_{b \in \pi} \left\{ \left[ \frac{\Gamma(t_b^j + 1)}{\prod_{k \in b} \Gamma(t_k^j + 1)} \right] \left[ \frac{\Gamma(\sum_{k \in b} \alpha_k)}{\prod_{k \in b} \Gamma(\alpha_k)} \right] \left[ \frac{\prod_{k \in b} \Gamma(\alpha_k + t_k^j)}{\Gamma(t_b^j + \sum_{k \in b} \alpha_k)} \right] \right\}$$

and

$$(3) \quad p(t_\pi^1, t_\pi^2 | A_\pi) = \left[ \frac{\Gamma(n_1 + 1) \Gamma(n_2 + 1)}{\prod_{b \in \pi} \Gamma(t_b^1 + 1) \Gamma(t_b^2 + 1)} \right] \left[ \frac{\Gamma(\sum_{b \in \pi} \beta_b)}{\prod_{b \in \pi} \Gamma(\beta_b)} \right] \left[ \frac{\prod_{b \in \pi} \Gamma(\beta_b + t_b^1 + t_b^2)}{\Gamma(n_1 + n_2 + \sum_{b \in \pi} \beta_b)} \right].$$

Let's look at some special cases to dissect this result.

Case 1. If  $\pi$  has a single block equal to the entire set of cell types  $\{1, 2, \dots, K\}$ , then  $t_b^j = n_j$  for both  $j = 1, 2$ , and the second formula reduces, correctly, to  $p(t_\pi^1, t_\pi^2 | A_\pi) = 1$ . Further,

$$p(t^j | t_\pi^j) = \left[ \frac{\Gamma(n_j + 1)}{\Gamma(n_1 + \sum_{k=1}^K \alpha_k)} \right] \left[ \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \right] \left[ \prod_{k=1}^K \frac{\Gamma(\alpha_k + t_k^j)}{\Gamma(t_k^j + 1)} \right]$$

which is the well-known Dirichlet-multinomial predictive distribution for counts  $t^j$  [15]. E.g, taking  $\alpha_k = 1$  for all types  $k$  we get the uniform distribution

$$p(t^j | t_\pi^j) = \frac{\Gamma(n_j + 1) \Gamma(K)}{\Gamma(n_j + K)}.$$

Case 2. At the opposite extreme,  $\pi$  has one block  $b$  for each class  $k$ . Then  $t_b^j = z_k^j$ , and  $p(t^j | t_\pi^j) = 1$ , and further, assuming  $\beta_b = \alpha_k$ ,

$$p(t_\pi^1, t_\pi^2 | A_\pi) = \left[ \frac{\Gamma(n_1 + 1) \Gamma(n_2 + 1)}{\prod_{k=1}^K \Gamma(t_k^1 + 1) \Gamma(t_k^2 + 1)} \right] \left[ \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \right] \left[ \frac{\prod_{k=1}^K \Gamma(\alpha_k + t_k^1 + t_k^2)}{\Gamma(n_1 + n_2 + \alpha_k)} \right].$$

which corresponds to Dirichlet-multinomial predictive distribution for counts  $t^1 + t^2$  since  $t^1$  and  $t^2$  are identical distributed in this case.

Regardless of the partition, log scale probabilities are readily evaluated given hyper-parameters  $\{\alpha_k\}$  and  $\{\beta_b\}$  and for cell-type counts  $t^1$  and  $t^2$ .

For asymptotic properties of the posterior probabilities, we demonstrated them in section 6.



## 4. SIMULATION STUDY

A simulation study was conducted to assess the performance of scDDboost in identifying DD genes. We simulate data by splatter[16] with approximate 200 cells each condition and 7 subtypes with proportions  $\phi$  and  $\psi$  from Figure 1 satisfying constraints:  $\phi_1 + \phi_2 = \psi_1 + \psi_2$ ,  $\phi_3 + \phi_4 + \phi_5 = \psi_3 + \psi_4 + \psi_5$  and  $\phi_6 + \phi_7 = \psi_6 + \psi_7$ . Each subtype has 10% genes to be differential expressed. We view the differences among subtypes by projecting transcripts profiles of cells into its first two principal components (figure 2). We observed subtypes are well separated, which is driven by genes with heterogeneous distribution between subtypes.

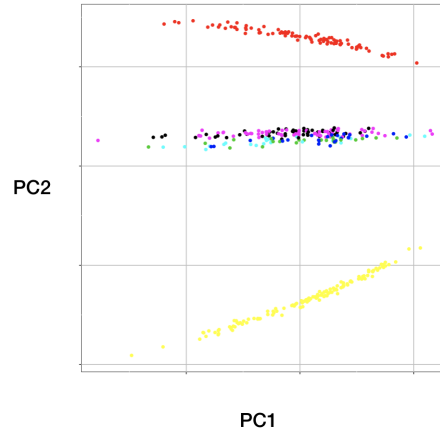


FIGURE 2. first two principal components of transcripts, which demonstrates the difference between subtypes. Even if we projected transcripts of cells into two dimensional space, we observe subtypes are well separated.

We determine the number of subtypes by searching a range of candidates (from 1 to 9 based on our empirical experience). Given number of subtypes, we obtain a subtype structure of cells, which will further be fed into computing the posterior probabilities. We visualize the change between posterior probabilities under number of clusters  $i$  and  $i + 1$  ( $i$  from 1 to 8). It typically remains stable when number of cluster is above a number that is smaller than 9 (figure 3). In the simulated data, the posterior probability become stable when we overestimate the number of subtypes. We found the true number of subtypes is 7 and correctly identify the subtypes of cells.

10% DE genes in each subtype results in total 8704 DD genes and 8669 ED genes in the mixture of the 7 subtypes. Table 1 are numbers of DD or DE genes identified by four methods (scDDboost, scDD, MAST and DESeq2) with target FDR at 5 %.

scDDboost identified most true DD genes, the reason is that mean expression shifts between conditions is not as significant as mean expression shifts between subtypes, which limits

	scDDboost	scDD	MAST	DESeq2
DD or DE genes	5126	1593	2559	3000
True positive	5094	1570	2508	2928
false positive	32	23	51	72

TABLE 1. number of true positive and false positive genes identified by four methods. Target FDR at 5%

the power of MAST and DESeq2. Our approach and scDD considered mixture structure underlying the transcripts but scDD did not use the whole genome information to infer mixture components, which leads to inaccurate clustering at gene level and reduce the power. scDDboost could correctly identify the subtypes of cells and thus are more sensitive to the mean expression change among subtypes. We also compare roc curves of scDDboost, scDD, MAST and DESeq2. (figure 4)

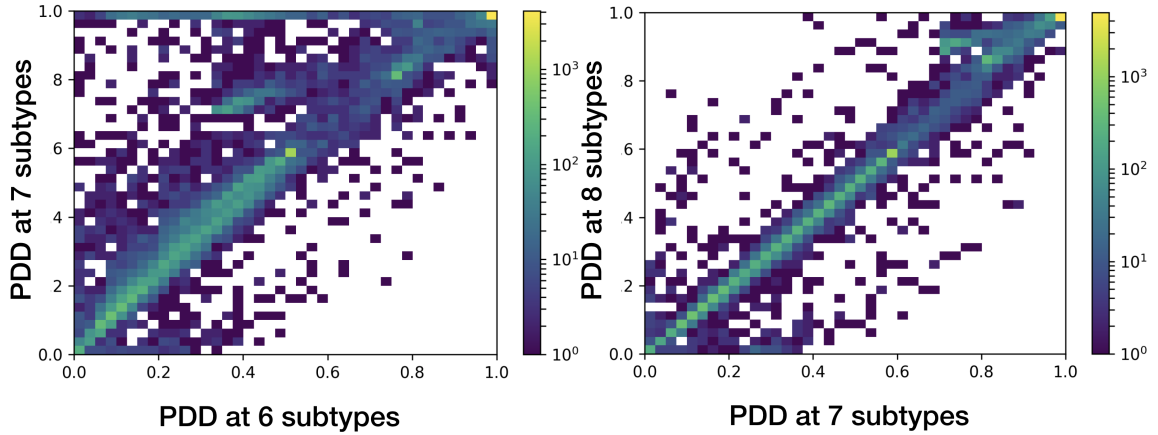


FIGURE 3. comparison of posterior probabilities of being DD among different number of subtypes, when we underestimate the number of subtypes, the difference is huge, see PDD between 6 subtypes and 7 subtypes. When we overestimate the number of subtypes, though inflating PDD but the variation of difference is small, from 6 to 8 subtypes the PDD become more linear related

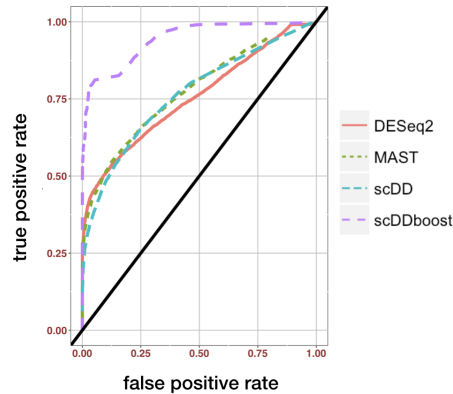


FIGURE 4. Roc curve of scDDboost, scDD, MAST and DESeq2, scDDboost has largest area under the roc curve. Roc curves of other three methods are similar. For those roc curve there is bigger difference at low level of false positive rate, as scDDboost identified twice many true DD genes as other methods.

Since we are modeling gene transcript within each subtype as negative binomial distributed and we only test one parameter(mean) change among subtypes. In some scenario, it could be insufficient to model the variability within subtype. Even though there is no mean expression change among subtypes but more subtle distributional change occurred among subtypes changed, EBSeq would fail to detect the discrepancies between subtypes, thus limit power of scDDboost.

## 5. EXAMPLES

We use ten datasets from conquer[9] to test performance of our method on real data. We compare our results with scDD[8], MAST[6] and DESeq2[7]

Data set	Compared cell subsets	Number of cells/condition	Organism	Ref
GSE45719	16-cell stage blastomere vs Mid blastocyst cell (92-94h post-fertilization)	50, 60	mouse	[17]
GSE45719null	16-cell stage blastomere	50	mouse	[17]
GSE48968-GPL13112	BMDC (1h LPS stimulation) vs BMDC(4h LPS stimulation)	96, 95	mouse	[18]
GSE48968-GPL13112null	BMDC (1h LPS stimulation)	96	mouse	[18]
GSE60749-GPL13112	v6.5 mouse embryonic stem cells, culture conditions: 2i+LIF vs v6.5 mouse embryonic stem cells, culture conditions: serum+LIF	90, 94	mouse	[19]
GSE60749-GPL13112null	v6.5 mouse embryonic stem cells, culture conditions: 2i+LIF	90	mouse	[19]
GSE74596	NKT0 vs NKT17	45,44	mouse	[20]
GSE74596null	NKT0	45	mouse	[20]
EMTAB2805	G1 vs G2M	96,96	mouse	[21]
EMTAB2805null	G1	96	mouse	[21]
GSE63818-GPL16791	Primordial Germ Cells, developmental stage: 7 week gestation vs Somatic Cells, developmental stage: 7 week gestation	39,26	mouse	[22]
GSE71585-GPL13112	Chrna2 tdTpositive vs Cux2 tdTpositive	84, 124	mouse	[23]
GSE71585-GPL13112null	Chrna2 tdTpositive	84	mouse	[23]
GSE75748	NPC vs DEC	64, 87	human	[24]
GSE75748	NPC	64	human	[24]
GSE75748	DEC vs EC	70, 64	human	[24]
GSE75748	DEC	70	human	[24]
GSE64016null	H1 exp1 vs H1 exp2	64, 87	human	[25]

TABLE 2. single cell transcripts profiles used for differential expression or distribution method evaluation

We have table of numbers of differentially expressed genes of each dataset by MAST and DESeq2, and numbers of differentially distributed genes of each dataset by scDDboost and scDD.

Data set	scDDboost	scDDboost-sc3	scDD	MAST	DESeq2	total number of genes
GSE45719	5758	4228	6416	5652	11202	45686
GSE48969-GPL13112	11691	9819	2080	3396	9542	45686
GSE60749-GPL13112	19215	19168	18074	13674	23178	45686
GSE74596	1942	1353	1099	540	3796	45686
EMTAB 2805	5295	3748	2202	1088	5391	45686
GSE63818-GPL16791	3948	3480	1365	873	8934	45686
GSE71585-GPL13112	2902	1460	1622	2572	7378	24057
NPC-DEC	4377	3211	5982	6666	8439	19037
DEC-EC	3402	3023	3818	5429	8127	19037
H1 exp1-H1 exp2	0	0	1300	2077	2841	16579

TABLE 3. number of genes detected as significantly DE or DD

We found that bulk method DESeq2 tends to have the most number of DE genes. But among single cell methods we found that scDDboost usually had the most number of DD genes. Further we observed quite a few genes uniquely identified by scDDboost are likely to have different distribution across conditions. For example, figure 5, we use violin plot to demonstrate the log expression profiles among DEC and EC.

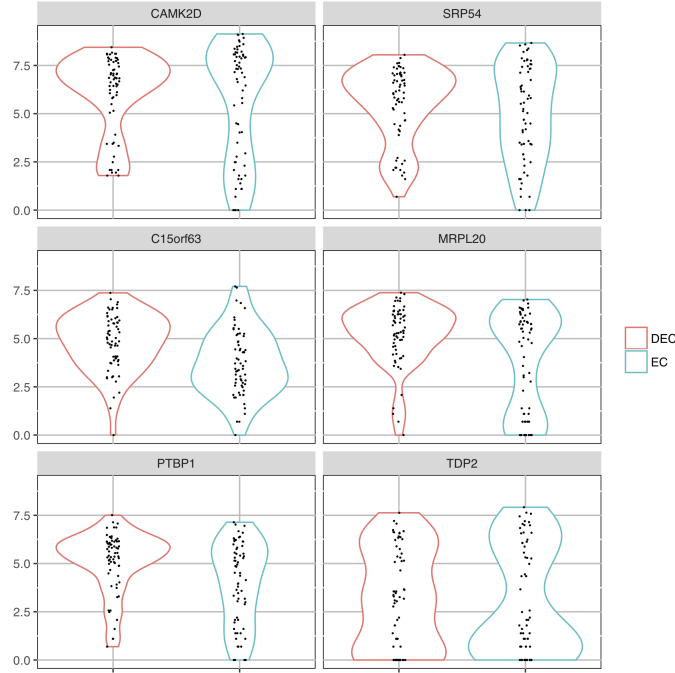


FIGURE 5. Densities of log transformed transcripts 6 DD genes uniquely identified by scDDboost, for data GSE75748, DEC vs. EC, We observe some of the genes are different distributed across conditions.

Although bulk methods seems to be the most powerful one, we found it also has a higher false discovery rate comparing to single cell methods. We validate false discovery rate on ten null datasets from table 1. For each null dataset, we randomly split the cells from one condition into two equal sized subsets and do DE analysis between those subsets. Since those two subsets of cells actually came from same condition, there should not be any differential distributed genes, any positive call would be a false positive. We repeat the random split and testing for five times on each null data set. We evaluate the type I error control for the methods returning nominal p-values, by recording the fraction of genes(with a valid p-value) that are assigned a nominal p-value below 0.05 (figure 5).

scDDboost could control FDR since we assume cells are sampled from population composed of different subtypes. Cells from one subtype are equal likely to be assigned to either one of the two subsets. Consequently, proportions of subtypes remain unchanged among the two subsets.

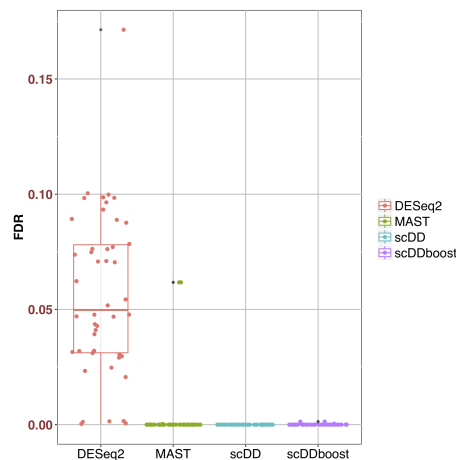


FIGURE 6. FDR of scDDboost, scDD, MAST and DESeq2 on null dataset from table 1, DESeq2 usually identify a lot but may lose the control of type I error. While other single cell methods could control FDR. \*\*\*This procedure for testing FDR, we randomly split a population into two samples, so it is highly likely proportion of subtypes remain same among these two samples, scDDboost always give small PDD and almost make no false positive call. This is kind different from really examining the FDR of scDDboost, where we have different proportions across conditions, but we do not want false call on those genes without mean expression change even there is proportion change. In this case, it actually reduce to the FDR of EBSeq whether we make correct posterior inference on DE pattern. The test of proportion is only done once, so scDDboost could still control FDR\*\*\*

D3E[26] is a distributional method that can identify bursting parameters of transcripts. Rate of promoter activation, rate of promoter inactivation and the rate of transcription when the promoter is in the active state are estimated by D3E. We investigate DD genes identified by scDDboost and their change of those three parameters on dataset EMTAB2805

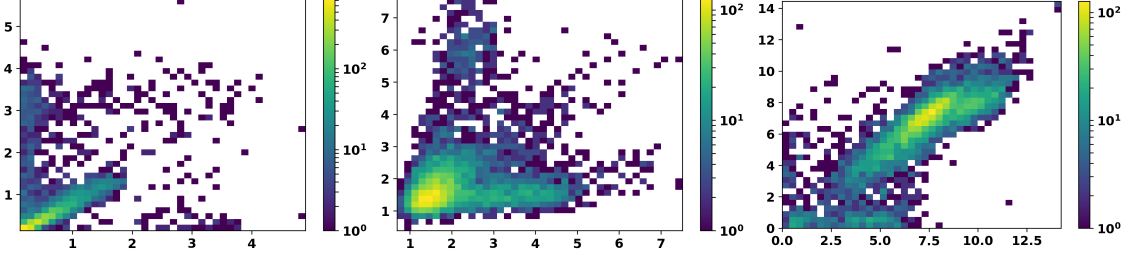


FIGURE 7. 2D histogram for bursting parameters of DD genes identified by scDDboost from dataset EMTAB2805 estimated by D3E. Left panel : comparison of rate of promoter activation between two conditions, similarly, middle panel : rate of promoter inactivation and right panel: rate of transcription when the promoter is in the active state. We observe that difference between transcription rate is smaller compare to difference between the activation and inactivation rate. \*\*\*other methods also observe similar phenomena, this is not unique to scDDboost, main reason is that estimation from D3E tends to give larger difference in activation and inactivation rate than transcription rate. We may argue the major factor to drive DD genes are activation and inactivation rate (proportions of different subtypes), so it make sense to consider mixture model like scDDboost.\*\*\*

We observed that DD genes identified by scDDboost tends to have similar transcription rate when the promoter is active across condition, while there are lots of variabilities in the action and inactivation rate. These results reveal that DD genes identified by scDDboost are driven by the change of activation and inactivation rates.

## 6. ROBUSTIFICATION

**6.1. bagging.** Cluster cells is an unsupervised learning, we do not know the true underlying partition and different number of clusters will lead to large differences in posterior probabilities of genes being differentially distributed (figure 7). We propose a modified bootstrap to stabilize our inferences. Instead of resample the cells, we resample the distance matrices of cells by adding noises to original distance matrix. Denote the original distance matrix as  $D = (d_{i,j})$ , for each time we random sample a vector  $e$  with length equal to number of cells and components are i.i.d. exponentially distributed. let  $w$  be the standard deviation of  $d_{i,j}$ . We resample a new  $\hat{D}$  by adding nosies:  $\hat{d}_{i,j} = d_{i,j} + e_i * w + e_j * w$ . For  $\hat{D}$  we still have triangle inequality held as  $\hat{d}_{i,j} + \hat{d}_{j,k} \geq \hat{d}_{i,k}$ , it is a valid distance matrix. For a fixed number of clusters, we average posterior probabilities over different distance matrices. And we select number of clusters  $K$  that posterior probabilities do not vary too



much under  $K$  and  $K + 1$ . From our empirical experience, it is typical  $K$  will not be larger than 8.

---

**Algorithm 2**


---

**Input:** cell by cell distance matrix  $D$  number of cluster  $K$  and parameter of noises  $\lambda$

**Output:** labelled partition of cells  $\hat{z}$

- 1: **procedure** RANDOM-C( $D, K, \lambda$ )
  - 2:   sampling noises  $e$  from exponential distribution with mean =  $\lambda$
  - 3:    $\hat{D}_{i,j} \leftarrow D_{i,j} + e_i + e_j$
  - 4:    $\hat{z} \leftarrow K$ -medoids on  $\hat{D}$
  - 5:   **return**  $\hat{z}$
- 

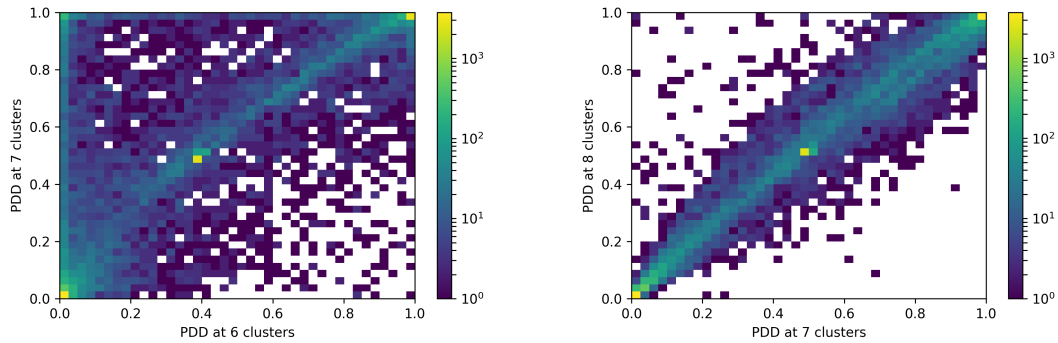


FIGURE 8. selecting number of subtypes for data GSE75748, we observe posterior probabilities become stable at more than 6 subtypes. Since increasing number of subtypes tends to decrease sample size of each subtypes, make complicate constraints for equivalent distribution and inflate estimated PDD. We select number of subtypes to be 7

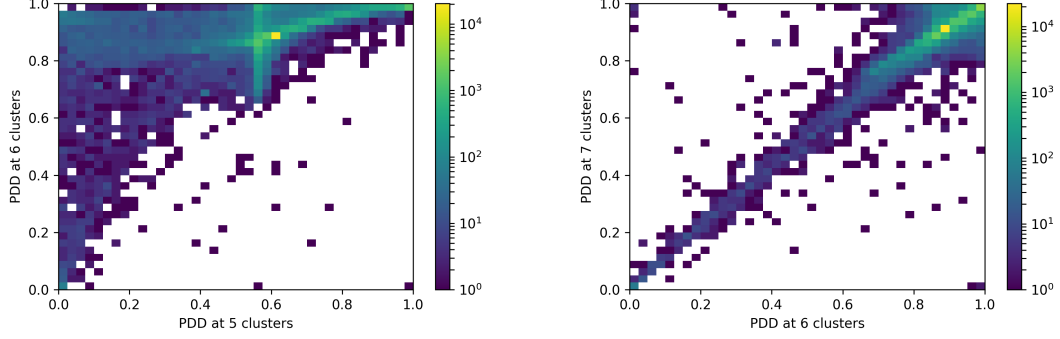


FIGURE 9. selecting number of subtypes for data GSE48968, we observe posterior probabilities become stable at more than 5 subtypes

---

### Algorithm 3

---

**Input:** Expression data matrix  $X$ , condition label vector  $y$ , number of clusters and iteration  $K, T$ , noise parameter  $\lambda$

**Output:** averaged posterior probabilities of genes being differential distributed

```

1: procedure SCDDBOOST1( $X, y, K, T, \lambda$ )
2:    $D \leftarrow$  cell by cell distance matrix  $X$ 
3:   for  $i = 0, i < T$  do
4:      $pt \leftarrow$  Random-c( $D, K$ )
5:      $P(\text{DD}_g|X, y) \leftarrow \text{PostI}(X, y, \hat{z})$ 
6:      $i \leftarrow i + 1$ 
7:   return averaged  $P(\text{DD}_g|X, y)$ 

```

---

## 7. ASYMPTOTIC PROPERTIES

To investigate asymptotic properties we first give the expression of posterior probability. Since there is no information favorable of any particular  $A_\pi$ , we select discrete uniform distribution as the prior for it, then the posterior probability is

$$(4) \quad p(A_\pi | t^1, t^2) = c * \sum_{\pi' \text{ refines } \pi} p(t^1 | t_{\pi'}^1) p(t^2 | t_{\pi'}^2) p(t_{\pi'}^1, t_{\pi'}^2 | A_{\pi'})$$

for a normalizing constant  $\frac{1}{c} = \sum_{\pi' \in \Pi} p(t^1 | t_{\pi'}^1) p(t^2 | t_{\pi'}^2) p(t_{\pi'}^1, t_{\pi'}^2 | A_{\pi'})$ .

Let  $\Omega = \{(\phi, \psi) : \sum_{i=1}^K \phi_i = \sum_{i=1}^K \psi_i = 1, \phi_i \geq 0, \psi_i \geq 0, i = 1, \dots, K\}$  be the whole space. There is a subset of  $\Omega$  we lack posterior inference. Let us first see an example:

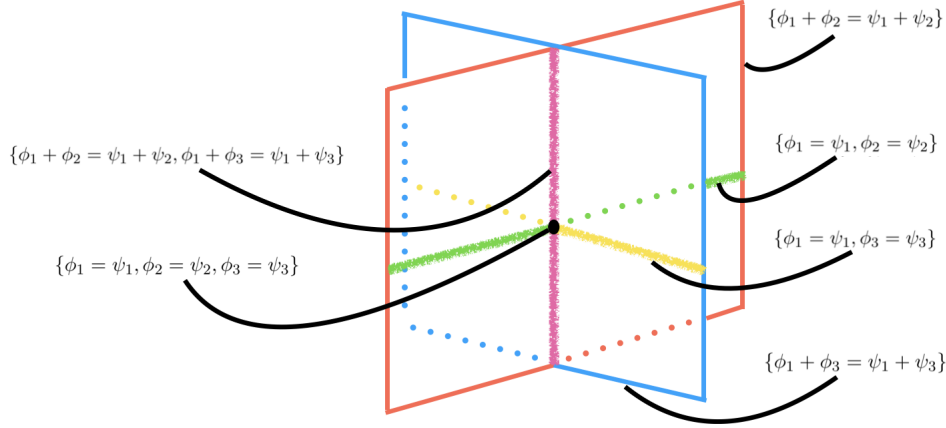


FIGURE 10. Four subtypes of cells, simplexes of  $(\phi, \psi)$  satisfying different constraints.

In figure 11, there are four subtypes, the rectangle with magenta boundary is a simplex  $A_{\pi_1} = \{(\phi, \psi) : \phi_1 + \phi_2 = \psi_1 + \psi_2\}$ , the rectangle with blue boundary is a simplex  $A_{\pi_2} = \{(\phi, \psi) : \phi_1 + \phi_3 = \psi_1 + \psi_3\}$ . The green line refers to  $A_{\pi_3} = \{(\phi, \psi) : \phi_1 = \psi_1, \phi_2 = \psi_2\}$ , the yellow line refers to  $A_{\pi_4} = \{(\phi, \psi) : \phi_1 = \psi_1, \phi_3 = \psi_3\}$ , the purple line refers to  $A_{\pi_5} = \{(\phi, \psi) : \phi_1 + \phi_2 = \psi_1 + \psi_2, \phi_1 + \phi_3 = \psi_1 + \psi_3\}$ , which is the intersection of  $A_{\pi_1}$  and  $A_{\pi_2}$ , and finally the black dot which is the intersection of those three lines refers to the simplex with finest partitions,  $\phi_i = \psi_i, \forall i = 1, \dots, 4$ . We lack posterior inference for  $(\phi, \psi)$  along the purple line except the black dot. While on the green line, yellow line and black dot, we have consistent posterior inference (theorem 2). To explain why some space lacking posterior inference and define such space, we define a special subset  $A_{\pi}^*$  of simplex  $A_{\pi}$ .  $A_{\pi}^* = A_{\pi} \setminus \bigcup_{\tilde{\pi} \text{ is not coarser than } \pi} A_{\tilde{\pi}}$ ,  $A_{\pi}^*$  is obtained by removing all intersection with other  $A_{\tilde{\pi}}$  (excluding those  $A_{\tilde{\pi}}$  that is superset of  $A_{\pi}$ ) from  $A_{\pi}$ . Since we removed those intersection parts. It is intuitive that  $A_{\pi}^*$  will be disjoint subsets of  $\Omega$ .

**Proposition 1.** *if  $\pi_1 \neq \pi_2$ , then  $A_{\pi_1}^* \cap A_{\pi_2}^* = \emptyset$*

Let  $Q = \Omega \setminus \bigcup_{\pi \in \Pi} A_{\pi}^*$ , and we have following proposition of the existence of  $Q$ .

**Proposition 2.** *Let  $K$  be number of subtypes. When  $K > 3, Q \neq \emptyset$ , when  $K \leq 3, Q = \emptyset$*

When number of subtypes bigger than three, we lack posterior inference on  $Q$ . To see that we can rewrite  $A_\pi^*$  as  $A_\pi^* = A_\pi \setminus \bigcup_{\tilde{\pi} \text{ is not coarser than } \pi} (A_{\tilde{\pi}} \cap A_\pi)$ ,  $\tilde{\pi}$  is not coarser than  $\pi$ , which is equivalently to say  $\pi$  is not refinement of  $\tilde{\pi}$ . By lemma 1,  $A_{\tilde{\pi}} \cap A_\pi$  is a lower dimensional subset of  $A_\pi$ . So  $A_\pi \setminus A_\pi^*$  is a lower dimensional subset of  $A_\pi$ . For posterior on  $Q$ , it degenerates to integral on a lower dimensional subset of the simplex associating with densities, which will vanish

**Proposition 3.** *When  $K > 3$ ,  $p(Q|z^1, z^2) = 0$*

But for  $(\phi, \psi) \in \Omega \setminus Q$ , we have consistent posterior inference. Assuming  $\alpha_i = 1, \forall i$  in (2) and  $\beta_b = \sum_{i \in b} \alpha_i$  in (3), plug in (4) then we have simplified

$$(5) \quad p(A_\pi | t^1, t^2) = \frac{1}{c'} \sum_{\pi' \in \text{RF}(\pi)} \prod_{b \in \pi'} \frac{\Gamma(\beta_b + t_b^1 + t_b^2)}{\Gamma(\beta_b + t_b^1) \Gamma(\beta_b + t_b^2)}$$

$c' = c / \frac{\Gamma(n+1)\Gamma(n+1)\Gamma(K)}{\Gamma(2n+K)}$  And we have theorem 3.

**Theorem 3.** *Let  $n = \min(n_1, n_2)$  be the smaller number of cells of two conditions and  $n_1 = O(n_2)$ , when parameter  $(\phi, \psi) \in \Omega \setminus Q$  we have*

$$p(A_\pi | t^1, t^2) \xrightarrow[n \rightarrow \infty]{a.s.} \begin{cases} 1 & \text{if } (\phi, \psi) \in A_\pi \\ 0 & \text{otherwise} \end{cases}$$

Things become more complicate when  $(\phi, \psi)$  falling into  $Q$ , we know  $p(Q|t^1, t^2)$  vanishes, but  $p(A_\pi | t^1, t^2)$  may not.

Recall  $N(\pi)$  represents number of blocks  $b$  in  $\pi$ . Let  $S = \{\pi, (\phi, \psi) \in A_\pi\}$ , which is the collection of partitions whose associated simplexes covering  $(\phi, \psi)$ . Let  $N^* = \max_{\pi \in S} N(\pi)$ , which is the max number of blocks of partitions from  $S$ . Let  $S^* = \{\pi, (\phi, \psi) \in A_\pi \text{ and } N(\pi) = N^*\}$ , which is the collection of partitions that covering  $(\phi, \psi)$  with number of blocks equal to the max number  $N^*$ .

For example, when  $K = 7$ , For a  $(\phi, \psi) \in A_{\pi_1} \cap A_{\pi_2} \cap A_{\pi_3}$ ,  $\pi_1 = \{\{1, 2, 3\}, \{4, 5, 6, 7\}\}$ ,  $\pi_2 = \{\{1, 6, 7\}, \{2, 4\}, \{3, 5\}\}$ ,  $\pi_3 = \{\{1, 2, 3, 4, 5, 6\}\}$ , and also  $(\phi, \psi)$  does not belong to any other simplex  $A_\pi$ . Then  $S = \{\pi_1, \pi_2, \pi_3\}$ ,  $N^* = 3$ ,  $S^* = \{\pi_2\}$ .

Denote components from right hand side of (5):  $\frac{1}{c'} \prod_{b \in \pi} \frac{\Gamma(\beta_b + t_b^1 + t_b^2)}{\Gamma(\beta_b + t_b^1) \Gamma(\beta_b + t_b^2)} = J(t^1, t^2, \pi)$ . We have theorem 4.

**Theorem 4.** *Following the setting in theorem 2, when parameter  $(\phi, \psi) \in Q$ , and we have*

$$J(t^1, t^2, \pi) \xrightarrow[n \rightarrow \infty]{a.s.} \begin{cases} m(\pi) & \pi \in S^* \\ 0 & \text{otherwise} \end{cases}$$

and  $\sum_{\pi \in S^*} m(\pi) = 1, m(\pi) > 0$

proofs are in the appendix.

Still using above example, in limiting case, we have  $p(A_{\pi_3}|t^1, t^2) = 1$ ,  $p(A_{\pi_2}|t^1, t^2) = 1$  and  $p(A_{\pi_1}|t^1, t^2) = 0$ . When the DE pattern is  $B_{\pi_1}$  for some genes. Since our underestimation of  $p(A_{\pi_1}|z^1, z^2) = 0$ , we will falsely classify those genes as differential distributed.

The asymptotic properties help us gain insight of the performance of our approach, scD-Dboost may work poorly, when  $(\phi, \psi) \in Q$ , we may underestimate the posterior probability of true proportion change pattern, which reduce the posterior probabilities of true negative and enlarge false positive rate.

## 8. DISCUSSION

subsectionextension to multi conditions

## APPENDIX A

**Lemma 1.** *If  $\pi_2$  is not refinement of  $\pi_1$  then  $A_{\pi_1} \cap A_{\pi_2}$  is a lower dimensional subset of  $A_{\pi_2}$*

Proof of theorem 2

*Proof.* by lemma 1, it is easy to verify.  $\square$

where  $p(t^1, t^2|\phi, \psi) = p(t^1|\phi)p(t^1|\psi)$ ,  $t^1|\phi \sim \text{multinomial}(n_1, \phi)$ ,  $t^2|\psi \sim \text{multinomial}(n_2, \psi)$ . Recall the definition of  $A_\pi = \{(\phi, \psi) : \Phi_b = \Psi_b\}$  and  $A_\pi$  is a simplex. Denote the finest partition as  $\pi_F = \{\{1\}, \{2\}, \dots, \{K\}\}$ , associated simplex  $A_{\pi_F} = \{(\phi, \psi) : \phi_i = \psi_i, i = 1, \dots, K\}$  for any two partition  $\pi_1$  and  $\pi_2$ , intersection of their associated simplex must not be empty since  $A_{\pi_F} \subset A_{\pi_1} \cap A_{\pi_2} \neq \emptyset$ . To discuss the issue of overlapping of simplex  $A_\pi$ , we first introduce some notations. The whole space  $\Omega = \{(\phi, \psi), \phi_i, \psi_i > 0 \text{ and } \sum_{i=1}^K \phi_i = \sum_{i=1}^K \psi_i = 1\}$  and we define the refinement and coarseness relationship between partitions, we say a partition  $\tilde{\pi}$  refines another partition  $\pi$  if  $\forall b \in \pi$  there exists  $s \subset \tilde{\pi}$  such that  $\cup_{b' \in s} b' = b$ . When  $\tilde{\pi}$  refines  $\pi$ , we say  $\tilde{\pi}$  is a refinement of (finer than)  $\pi$  or  $\pi$  is a coarseness of (coarser than)  $\tilde{\pi}$ . Observe that if  $\pi'$  refines  $\pi$ , then  $A_\pi \cap A_{\pi'} = A_{\pi'}$ ,  $\int_{A_\pi \cap A_{\pi'}} p(z^1, z^2|\phi, \psi)p(\phi, \psi|A_{\pi'})d\phi d\psi = \int_{A_{\pi'}} p(t^1, t^2|\phi, \psi)p(\phi, \psi|A_{\pi'})d\phi d\psi$ . When  $\pi'$  is not refinement of  $\pi$ , we need to know the dimension of  $A_\pi \cap A_{\pi'}$ . Consider a map  $f : b \rightarrow v$ , which maps the block  $b$  to a vector  $v \in \{0, 1\}^K$ , the  $i$ th component of  $v$  is  $1_{\{i \in b\}}$ . And denote  $\dim(S)$  be the dimension of space  $S$ .  $A_\pi$  can be equivalently defined as  $A_\pi = \{(\phi, \psi) : M_\pi * (\phi - \psi) = 0\}$ ,  $M_\pi$  is a matrix with rows be  $v_b = f(b)$ ,  $\forall b \in \pi$ , that is to say  $(\phi, \psi)$  are in the null space of linear transformation  $M_\pi$ . We have following lemma

Proof of lemma 1

*Proof.* Let  $V$  denote the orthogonal space of  $\phi - \psi$ , when  $(\phi, \psi) \in A_{\pi_1} \cap A_{\pi_2}$ , and  $\dim(A_{\pi_1} \cap A_{\pi_2}) = \dim(\phi - \psi) + \dim(\psi) = 2K - \dim(V) - 1$ . Also let  $\pi_1 = \{b_1^1, \dots, b_s^1\}$ ,  $\pi_2 = \{b_1^2, \dots, b_t^2\}$ . The corresponding vectors are  $v_1^1, \dots, v_s^1$  and  $v_1^2, \dots, v_t^2$ . We claim there must be a  $b_i^1 \in \pi$  whose corresponding  $v_i^1$  is linear independent with  $v_1^2, \dots, v_t^2$ . If not, for every  $v_i^1$  there exists  $\alpha_1^i, \dots, \alpha_t^i$  such that

$$v_i^1 = \sum_{j=1}^t \alpha_j^i v_j^2 \quad (*)$$

If  $b_j^2 \cap b_i^1 \neq \emptyset$ , then multiply  $v_j^2$  on both sides of (\*), we obtain  $v_i^1 * v_j^2 = \alpha_j^i (v_j^2)^2$ , as  $v_j^2$  are orthogonal vectors, and  $v_i^1 * v_j^2 > 0$  implies  $\alpha_j^i > 0$ . Consider  $x = f(b_j^2 \setminus b_i^1)$ , we have  $x * v_i^1 = 0$  and we multiply  $x$  on both sides of (\*) to obtain  $\alpha_j^i v_j^2 * x = 0$ , thus  $x$  must be zero vector and  $b_j^2 \setminus b_i^1 = \emptyset$ , which implies  $b_j^2 \subset b_i^1$ . That is to say when  $b_j^2 \cap b_i^1 \neq \emptyset$ ,  $b_j^2$  must be subset of  $b_i^1$ . So  $b_i^1$  is union of some blocks in  $\pi_2$ . Which implies  $\pi_2$  is refinement of  $\pi_1$ ,

contradiction.

Consequently there exists  $b \in \pi_1$  with  $v(b)$  linear independent with  $v(b'), b' \in \pi_2$ .  $\dim(V)$  is at least  $N(\pi_2) + 1$ ,  $\dim(A_{\pi_1} \cap A_{\pi_2}) < \dim(A_{\pi_2})$   $\square$

Proof of theorem 3 and theorem 4

*Proof.* Given the condition that  $\alpha_k = 1, \forall k$  and  $\beta_b = \sum_{k \in b} \alpha_k$ , recall  $p(A_\pi | t^1, t^2) = \sum_{\pi' \in \text{RF}(\pi)} J(t^1, t^2, \pi')$  and  $J(t^1, t^2, \pi) = \frac{1}{c'} \prod_{b \in \pi} \frac{\Gamma(\beta_b + t_b^1 + t_b^2)}{\Gamma(\beta_b + t_b^1) \Gamma(\beta_b + t_b^2)}$ . Assuming there are  $K$  subgroups, since  $n_1$  and  $n_2$  goes to infinite at same rate, for simplicity we assume  $n = \sum_{i=1}^K t_i^1 = \sum_{i=1}^K t_i^2$ ,  $t^1 \sim \text{multinomial}(\phi)$ ,  $t^2 \sim \text{multinomial}(\psi)$  and  $t_b^1 = \sum_{i \in b} z_i^1$  and  $t_b^2 = \sum_{i \in b} z_i^2$ , so  $t_b^1 \sim \text{binomial}(n, \Phi_b)$  and  $t_b^2 \sim \text{binomial}(n, \Psi_b)$ , where  $\Phi_b = \sum_{i \in b} \phi_i$  and  $\Psi_b = \sum_{i \in b} \psi_i$ . Let  $f(n, b) = \frac{\Gamma(\beta_b + t_b^1 + t_b^2)}{\Gamma(\beta_b + t_b^1) \Gamma(\beta_b + t_b^2)}$ , then

$$J(z^1, z^2, \pi) \propto \prod_{b \in \pi} f(n, b)$$

$\log f(n, b) = \log(\Gamma(\beta_b + t_b^1 + t_b^2)) - \log(\Gamma(\beta_b + t_b^1)) - \log(\Gamma(\beta_b + t_b^2))$ , notice that  $t_b^1, t_b^2$  and  $\beta_b$  are integers, and when  $x$  is integer,  $\Gamma(x)$  is the factorial of  $(x-1)$ . We have  $\log f(n, b) = \log((\beta_b + t_b^1 + t_b^2 - 1)!) - \log((\beta_b + t_b^1 - 1)!) - \log((\beta_b + t_b^2 - 1)!) = \log((\beta_b + t_b^1 + t_b^2 - 1)!) - \log((\beta_b + t_b^1 - 1)!) - \log((\beta_b + t_b^2 - 1)!) \approx (\beta_b + t_b^1 + t_b^2 - 1) \log(\beta_b + t_b^1 + t_b^2 - 1) - (\beta_b + t_b^1 - 1) \log(\beta_b + t_b^1 - 1) - (\beta_b + t_b^2 - 1) \log(\beta_b + t_b^2 - 1) + O(\log(n))$ .

Plug into  $f(n, b)$  we have:

$$\log f(n, b) \approx (\beta_b + t_b^1 - 1) \log(1 + \frac{t_b^2}{\beta_b + t_b^1 - 1}) + (\beta_b + t_b^2 - 1) \log(1 + \frac{t_b^1}{\beta_b + t_b^2 - 1}) + O(\log(n))$$

as  $\beta_b \log(\beta_b + t_b^1 + t_b^2 - 1) \sim O(\log(n))$  and by law of large number and slusky's theorem,  $\log(1 + \frac{t_b^2}{\beta_b + t_b^1 - 1}) \rightarrow \log(1 + \frac{\Psi_b}{\Phi_b})$ ,  $\log(1 + \frac{t_b^1}{\beta_b + t_b^2 - 1}) \rightarrow \log(1 + \frac{\Phi_b}{\Psi_b})$  a.s. and  $\frac{\log f(n, b)}{n} \rightarrow \Phi_b \log(1 + \frac{\Psi_b}{\Phi_b}) + \Psi_b \log(1 + \frac{\Phi_b}{\Psi_b})$  a.s. We have:

$$\frac{\log(\prod_{b \in \pi} f(n, b))}{n} \rightarrow \sum_b [\Phi_b \log(1 + \frac{\Psi_b}{\Phi_b}) + \Psi_b \log(1 + \frac{\Phi_b}{\Psi_b})] \quad a.s.$$

To find the maxima  $(\Phi, \Psi)$ , we fix  $\Psi$  and let  $C = \frac{\log(\prod_{b \in \pi} f(n, b))}{n} + \lambda(\sum_{b \in \pi} \Phi_b - 1)$ , we have

$\frac{\partial C}{\partial \Phi_b} = \log(1 + \frac{\Psi_b}{\Phi_b}) + \lambda$ , stationary point is  $\Phi_b = \Psi_b, \forall b$ . and for the hessian matrix  $\frac{\partial^2 C}{\partial \Phi_b^2} = -\frac{\Psi_b}{\Phi_b^2 + \Phi_b \Psi_b} < 0$  and  $\frac{\partial^2 C}{\partial \Phi_b \partial \Phi_{b'}} = 0$ , if  $b \neq b'$ , that is to say the hessian matrix is a diagonal matrix with every diagonal elements to be negative, so it is negative definite, and our objective function is concave. The maxima is the stationary point  $\Phi = \Psi$ . And when  $\Phi = \Psi$ ,  $\frac{\log(\prod_{b \in \pi} f(n, b))}{n} = 2 \ln(2)$  a constant not dependent on partition  $\pi$  and  $\Phi$ . That is to

say if  $(\phi, \psi) \in A_{\pi_1} \cap A_{\pi_2}$  and  $(\phi, \psi) \notin A_{\pi_3}$ . Then we would have  $\lim_{n \rightarrow \infty} \frac{\log(\prod_{b \in \pi_1} f(n, b))}{n} = \lim_{n \rightarrow \infty} \frac{\log(\prod_{b \in \pi_2} f(n, b))}{n}$  and  $\lim_{n \rightarrow \infty} [\frac{\log(\prod_{b \in \pi_1} f(n, b))}{n} - \frac{\log(\prod_{b \in \pi_3} f(n, b))}{n}] = c > 0$ , which implies:

$$(A) \quad \frac{J(t^1, t^2, \pi_3)}{J(t^1, t^2, \pi_1)} \rightarrow 0 \quad a.s.$$

To investigate the limit of  $\frac{J(t^1, t^2, \pi_1)}{J(t^1, t^2, \pi_2)}$ , We use inequalities that  $\sqrt{2\pi}n^{n+\frac{1}{2}}e^{-n} \leq n! \leq en^{n+\frac{1}{2}}e^{-n}$  holds for all nonnegative integers  $n$ . Plug in  $f(n, b)$ , we have:

$$(1) \quad \beta_b + \log\sqrt{2\pi} - 3 + g(n, b) \leq f(n, b) \leq \beta_b - 2\log\sqrt{2\pi} + g(n, b)$$

$$g(n, b) = (\beta_b + t_b^1 - \frac{1}{2})\log(1 + \frac{t_b^2}{\beta_b + t_b^1 - 1}) + (\beta_b + t_b^2 - \frac{1}{2})\log(1 + \frac{t_b^1}{\beta_b + t_b^2 - 1}) - (\beta_b - \frac{1}{2})\log(\beta_b + t_b^1 + t_b^2 - 1)$$

Based on inequalities (1),  $\sum_{b \in \pi} f(n, b)$  only differ with  $\sum_{b \in \pi} g(n, b)$  by a constant. By Taylor's expansion  $\log(1 + x) = \log 2 + \frac{1}{2}(x - 1) + O((x - 1)^2)$ , we have  $\log(1 + \frac{t_b^2}{\beta_b + t_b^1 - 1}) = \log 2 + \frac{1}{2}(\frac{t_b^1 - t_b^2 + 1 - \beta_b}{\beta_b + t_b^1 - 1}) + O_p((\frac{t_b^1 - t_b^2 + 1 - \beta_b}{\beta_b + t_b^1 - 1})^2)$  and under condition  $\Phi_b = \Psi_b$ ,  $\frac{(t_b^1 - t_b^2 + 1 - \beta_b)^2}{\beta_b + t_b^1 - 1}$  is  $O_p(1)$ . Plug in  $g(n, b)$

$$g(n, b) = \log 2 * t_b^1 + \log 2 * t_b^2 - (\beta_b - \frac{1}{2})\log(\beta_b + t_b^1 + t_b^2 - 1) + O_p(1)$$

and sum up

$$(2) \quad \sum_{b \in \pi} g(n, b) = 2n\log 2 - \sum_{b \in \pi} (\beta_b - \frac{1}{2})\log(\beta_b + t_b^1 + t_b^2 - 1) + O_p(1)$$

Notice that when two partition  $\pi_1, \pi_2$  have same number of blocks  $b$  and  $\Phi_b = \Psi_b, \forall b \in \pi_1 \cup \pi_2$ ,

$$\begin{aligned} \sum_{b \in \pi_1} g(n, b) - \sum_{b' \in \pi_2} g(n, b') &= \sum_{b' \in \pi_2} (\beta_{b'} - \frac{1}{2})\log(\beta_{b'} + t_{b'}^1 + t_{b'}^2 - 1) - \sum_{b \in \pi_1} (\beta_b - \frac{1}{2})\log(\beta_b + t_b^1 + t_b^2 - 1) + O_p(1) \\ &= \sum_{b' \in \pi_2} (\beta_{b'} - \frac{1}{2})\log(\frac{\beta_{b'} + t_{b'}^1 + t_{b'}^2 - 1}{n}) - \sum_{b \in \pi_1} (\beta_b - \frac{1}{2})\log(\frac{\beta_b + t_b^1 + t_b^2 - 1}{n}) \\ &\quad + \sum_{b' \in \pi_2 - \frac{1}{2}} (\beta_{b'} - \frac{1}{2})\log(n) - \sum_{b \in \pi_1 - \frac{1}{2}} (\beta_b - \frac{1}{2})\log(n) + O_p(1) \\ &= O_p(1) + \sum_{b \in \pi_1} \frac{1}{2}\log(n) - \sum_{b' \in \pi_2} \frac{1}{2}\log(n) \\ &= O_p(1) \end{aligned}$$



When  $\pi_1$  and  $\pi_2$  have same number of blocks,

$$(B) \quad \frac{J(t^1, t^2, \pi_1)}{J(t^1, t^2, \pi_2)} \rightarrow O_p(1) \quad a.s.$$

When  $\pi_1$  have less blocks than  $\pi_2$ ,  $\sum_{b' \in \pi_2} g(n, b') - \sum_{b \in \pi_1} g(n, b) = O_p(\log(n))$

$$(C) \quad \frac{J(t^1, t^2, \pi_1)}{J(t^1, t^2, \pi_2)} \rightarrow 0 \quad a.s.$$

□

## REFERENCES

- [1] T. Nawy, “Single-cell sequencing,” *Nature Methods*, vol. 11, pp. 18 EP –, 12 2013. [Online]. Available: <http://dx.doi.org/10.1038/nmeth.2771>
- [2] R. Bacher and C. Kendzierski, “Design and computational analysis of single-cell rna-sequencing experiments,” *Genome Biology*, vol. 17, no. 1, p. 63, 2016. [Online]. Available: <https://doi.org/10.1186/s13059-016-0927-y>
- [3] V. Y. Kiselev, K. Kirschner, M. T. Schaub, T. Andrews, A. Yiu, T. Chandra, K. N. Natarajan, W. Reik, M. Barahona, A. R. Green, and M. Hemberg, “Sc3: consensus clustering of single-cell rna-seq data,” *Nature Methods*, vol. 14, pp. 483 EP –, 03 2017. [Online]. Available: <http://dx.doi.org/10.1038/nmeth.4236>
- [4] P. Lin, M. Troup, and J. W. K. Ho, “Cidr: Ultrafast and accurate clustering through imputation for single-cell rna-seq data,” *Genome Biology*, vol. 18, no. 1, p. 59, 2017. [Online]. Available: <https://doi.org/10.1186/s13059-017-1188-0>
- [5] E. Pierson and C. Yau, “Zifa: Dimensionality reduction for zero-inflated single-cell gene expression analysis,” *Genome Biology*, vol. 16, no. 1, p. 241, 2015. [Online]. Available: <https://doi.org/10.1186/s13059-015-0805-z>
- [6] G. Finak, A. McDavid, M. Yajima, J. Deng, V. Gersuk, A. K. Shalek, C. K. Slichter, H. W. Miller, M. J. McElrath, M. Prlic, P. S. Linsley, and R. Gottardo, “Mast: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell rna sequencing data,” *Genome Biology*, vol. 16, no. 1, p. 278, 2015. [Online]. Available: <https://doi.org/10.1186/s13059-015-0844-5>
- [7] M. I. Love, W. Huber, and S. Anders, “Moderated estimation of fold change and dispersion for rna-seq data with deseq2,” *Genome Biology*, vol. 15, no. 12, p. 550, 2014. [Online]. Available: <https://doi.org/10.1186/s13059-014-0550-8>
- [8] K. D. Korthauer, L.-F. Chu, M. A. Newton, Y. Li, J. Thomson, R. Stewart, and C. Kendzierski, “A statistical approach for identifying differential distributions in single-cell rna-seq experiments,” *Genome Biology*, vol. 17, no. 1, p. 222, 2016. [Online]. Available: <https://doi.org/10.1186/s13059-016-1077-y>
- [9] C. Sonesson and M. D. Robinson, “Bias, robustness and scalability in differential expression analysis of single-cell rna-seq data,” *bioRxiv*, 2017. [Online]. Available: <https://www.biorxiv.org/content/early/2017/05/28/143289>

- [10] N. Leng, J. A. Dawson, J. A. Thomson, V. Ruotti, A. I. Rissman, B. M. G. Smits, J. D. Haag, M. N. Gould, R. M. Stewart, and C. Kendziorski, “Ebseq: an empirical bayes hierarchical model for inference in rna-seq experiments,” *Bioinformatics*, vol. 29, no. 8, pp. 1035–1043, 2013. [Online]. Available: + <http://dx.doi.org/10.1093/bioinformatics/btt087>
- [11] R. Bacher, L.-F. Chu, N. Leng, A. P. Gasch, J. A. Thomson, R. M. Stewart, M. Newton, and C. Kendziorski, “Scnorm: robust normalization of single-cell rna-seq data,” *Nature Methods*, vol. 14, pp. 584 EP –, 04 2017. [Online]. Available: <http://dx.doi.org/10.1038/nmeth.4263>
- [12] D. B. Dahl, “Modal clustering in a class of product partition models,” *Bayesian Anal.*, vol. 4, no. 2, pp. 243–264, 06 2009. [Online]. Available: <https://doi.org/10.1214/09-BA409>
- [13] A. Strehl and J. Ghosh, “Cluster ensembles — a knowledge reuse framework for combining multiple partitions,” *J. Mach. Learn. Res.*, vol. 3, pp. 583–617, Mar. 2003. [Online]. Available: <https://doi.org/10.1162/153244303321897735>
- [14] B. Dickey J., Lientz, “The weighted likelihood ratio, sharp hypotheses, and the order of a markov chain.” *Ann. Math. Statist.*, vol. 41, no. 1, p. 214, 1970. [Online]. Available: <https://projecteuclid.org/euclid.aoms/1177697203>
- [15] U. Wagner and A. Taudes, “A multivariate polya model of brand choice and purchase incidence,” *Marketing Science*, vol. 5, no. 3, pp. 219–244, Aug. 1986. [Online]. Available: <http://dx.doi.org/10.1287/mksc.5.3.219>
- [16] L. Zappia, B. Phipson, and A. Oshlack, “Splatter: simulation of single-cell rna sequencing data,” *Genome Biology*, vol. 18, no. 1, p. 174, 2017. [Online]. Available: <https://doi.org/10.1186/s13059-017-1305-0>
- [17] Q. Deng, D. Ramsköld, B. Reinius, and R. Sandberg, “Single-cell rna-seq reveals dynamic, random monoallelic gene expression in mammalian cells,” *Science*, vol. 343, no. 6167, pp. 193–196, 2014. [Online]. Available: <http://science.sciencemag.org/content/343/6167/193>
- [18] A. K. Shalek, R. Satija, J. Shuga, J. J. Trombetta, D. Gennert, D. Lu, P. Chen, R. S. Gertner, J. T. Gaublot, N. Yosef, S. Schwartz, B. Fowler, S. Weaver, J. Wang, X. Wang, R. Ding, R. Raychowdhury, N. Friedman, N. Hacohen, H. Park, A. P. May, and A. Regev, “Single-cell rna-seq reveals dynamic paracrine control of cellular variation,” *Nature*, vol. 510, pp. 363 EP –, 06 2014. [Online]. Available: <http://dx.doi.org/10.1038/nature13437>
- [19] R. M. Kumar, P. Cahan, A. K. Shalek, R. Satija, A. Jay DaleyKeyser, H. Li, J. Zhang, K. Pardee, D. Gennert, J. J. Trombetta, T. C. Ferrante, A. Regev, G. Q. Daley, and J. J. Collins, “Deconstructing transcriptional heterogeneity in pluripotent stem cells,” *Nature*, vol. 516, pp. 56 EP –, 12 2014. [Online]. Available: <http://dx.doi.org/10.1038/nature13920>
- [20] I. Engel, G. Seumois, L. Chavez, D. Samaniego-Castruita, B. White, A. Chawla, D. Mock, P. Vijayanand, and M. Kronenberg, “Innate-like functions of natural killer t cell subsets result from highly divergent gene programs,” *Nature Immunology*, vol. 17, pp. 728 EP –, 04 2016. [Online]. Available: <http://dx.doi.org/10.1038/ni.3437>

- [21] F. Buettner, K. N. Natarajan, F. P. Casale, V. Proserpio, A. Scialdone, F. J. Theis, S. A. Teichmann, J. C. Marioni, and O. Stegle, “Computational analysis of cell-to-cell heterogeneity in single-cell rna-sequencing data reveals hidden subpopulations of cells,” *Nature Biotechnology*, vol. 33, pp. 155 EP –, 01 2015. [Online]. Available: <http://dx.doi.org/10.1038/nbt.3102>
- [22] F. Guo, L. Yan, H. Guo, L. Li, B. Hu, Y. Zhao, J. Yong, Y. Hu, X. Wang, Y. Wei, W. Wang, R. Li, J. Yan, X. Zhi, Y. Zhang, H. Jin, W. Zhang, Y. Hou, P. Zhu, J. Li, L. Zhang, S. Liu, Y. Ren, X. Zhu, L. Wen, Y. Q. Gao, F. Tang, and J. Qiao, “The transcriptome and dna methylome landscapes of human primordial germ cells,” *Cell*, vol. 161, no. 6, pp. 1437–1452, 2017/12/05. [Online]. Available: <http://dx.doi.org/10.1016/j.cell.2015.05.015>
- [23] B. Tasic, V. Menon, T. N. Nguyen, T. K. Kim, T. Jarsky, Z. Yao, B. Levi, L. T. Gray, S. A. Sorensen, T. Dolbeare, D. Bertagnolli, J. Goldy, N. Shapovalova, S. Parry, C. Lee, K. Smith, A. Bernard, L. Madisen, S. M. Sunkin, M. Hawrylycz, C. Koch, and H. Zeng, “Adult mouse cortical cell taxonomy revealed by single cell transcriptomics,” *Nature Neuroscience*, vol. 19, pp. 335 EP –, 01 2016. [Online]. Available: <http://dx.doi.org/10.1038/nn.4216>
- [24] L.-F. Chu, N. Leng, J. Zhang, Z. Hou, D. Mamott, D. T. Vereide, J. Choi, C. Kendzierski, R. Stewart, and J. A. Thomson, “Single-cell rna-seq reveals novel regulators of human embryonic stem cell differentiation to definitive endoderm,” *Genome Biology*, vol. 17, no. 1, p. 173, 2016. [Online]. Available: <https://doi.org/10.1186/s13059-016-1033-x>
- [25] N. Leng, L.-F. Chu, C. Barry, Y. Li, J. Choi, X. Li, P. Jiang, R. M. Stewart, J. A. Thomson, and C. Kendzierski, “Oscope identifies oscillatory genes in unsynchronized single-cell rna-seq experiments,” *Nature Methods*, vol. 12, pp. 947 EP –, 08 2015. [Online]. Available: <http://dx.doi.org/10.1038/nmeth.3549>
- [26] M. Delmans and M. Hemberg, “Discrete distributional differential expression (d3e) - a tool for gene expression analysis of single-cell rna-seq data,” *BMC Bioinformatics*, vol. 17, no. 1, p. 110, 2016. [Online]. Available: <https://doi.org/10.1186/s12859-016-0944-6>

*E-mail address:* newton@biostat.wisc.edu