

## Supplementary Material

Main document: *A compositional model to assess expression changes from single-cell RNA-seq data*

Authors: Ma, Korthauer, Kendzierski, and Newton

Version: May 5, 2019

This supplement is organized to match the sectioning of the main document. In summary,

1. Introduction
  - R package
2. Modeling
  - Data Structure, Sampling Model, and Parameters
    - Proof of Theorem 2
  - Method Structure and Clustering
    - EBSeq
    - modalClust
    - Randomized  $K$ –means
    - Selecting  $K$
  - Double Dirichlet Mixture
    - Proof of Properties 1-6 and Theorem 3
3. Numerical Experiments
  - Synthetic data, splatter
  - Empirical study, conquer and Null case
  - Robustness
4. Posterior consistency
  - proof of theorem 4

## 1. Introduction.

1.1. *R package*. Reference can be found at github site ...

\*\*on scDDboost, web page, etc\*\*

## 2. Modeling.

2.1. *Data Structure, Sampling Model, and Parameters*. Proof of Theorem 2.

PROOF. Recall  $\theta = (\phi, \psi, \mu, \sigma)$ . Through the sampling procedure of our model (Figure 3), assuming we known number of cells within each conditions  $(n_1, n_2)$ . We have  $z^1, z^2$  are multinomial draw given  $\phi$  and  $\psi$ , thus the generation of  $y, z$  only depends on  $(\phi, \psi)$  Also given  $z$ ,  $X_{g,c}$  is sampled through  $NB(\mu_{g,z_c}, \sigma_g)$ , only depends on  $(\mu, \sigma)$ . Thus  $P(X, y, z|\theta) = P(y, z|\phi, \psi)P(X|z, \mu, \sigma)$ , and we independently give priors for  $(\mu, \sigma)$  and  $(\phi, \psi)$  By the Baye's rule,

$$\begin{aligned} P(\theta|X, y, z) &\propto P(X, y, z|\theta)P(\theta) \\ P(X, y, z|\theta)P(\theta) &= P(y, z|\phi, \psi)P(X|z, \mu, \sigma)P(\mu, \sigma|z)P(\phi, \psi) \\ P(\phi, \psi|y, z) &\propto P(y, z|\phi, \psi)P(\phi, \psi) \\ P(\mu, \sigma|X, z) &\propto P(X|z, \mu, \sigma)P(\mu, \sigma|z) \\ \text{Thus } P(\theta|X, y, z) &\propto P(\phi, \psi|y, z)P(\mu, \sigma|X, z) \end{aligned}$$

From this we know

1. Given  $X, y$  and  $z$ ,  $(\phi, \psi) \perp (\mu, \sigma)$
2. Given condition and subtypes label  $y, z$ ,  $(\phi, \psi)$  is independent with  $X$
3. Given  $X$  and  $z$ ,  $(\mu, \sigma)$  is independent with  $y$

Thus we have  $P(A_\pi \cap M_{g,\pi}|X, y, z) = P(A_\pi|y, z) P(M_{g,\pi}|X, z)$ . □

2.2. *Method Structure and Clustering*.

2.2.1. *EBSeq*. In this subsection, we go through how we implement and modified EBSeq to get  $P(M_{g,\pi}|X, z)$ .

Suppose we have  $K$  subtypes, let  $X_g^I = X_{g,1}^I, \dots, X_{g,S_1}^I$  denote transcripts at gene  $g$  from subtype  $I, I = 1, \dots, K$ . In the EBSeq, it assumed that counts within subtype  $I$  are distributed as Negative Binomial:  $X_{g,s}^I | r_{g,s}, q_g^I \sim NB(r_{g,s}, q_g^I)$ . Due to sample specific sizefactor in the raw counts,  $r$  is made sample specific. However, we are dealing with normalized counts rather than raw counts in EBSeq, We make  $r$  shared at gene level across all samples, i.e.  $X_{g,s}^I | \sigma_g, q_g^I \sim NB(\sigma_g, q_g^I)$

$$P(X_{g,s}^I | \sigma_g, q_g^I) = \binom{X_{g,s}^I + \sigma_g - 1}{X_{g,s}^I} (1 - q_g^I)^{X_{g,s}^I} (q_g^I)^{\sigma_g}$$

and  $\mu_{g,s}^I = \sigma_g(1 - q_g^I)/q_g^I$ ; For the ease of later deriving the density kernel  $f$ , we use  $q$  rather than  $\mu$  to parameterize the NB.

From EBSeq, we assumed a prior distribution on  $q_g^I : q_g^I | \alpha, \beta^g \sim \text{Beta}(\alpha, \beta^g)$ . The hyperparameter  $\alpha$  is shared by whole genome and  $\beta^g$  is gene specific.

We force the size factor to be 1 for all cells and use the same procedure from EBSeq to estimate number of failure parameter  $\sigma_g$ . Namely, we have

1. gene-level sample mean  $m_g = \frac{1}{n} \sum_{s=1}^n X_{g,s}$ , where  $n = n_1 + n_2$  is the total number of cells

2. averaged sample variance within subtype  $v_g = \frac{1}{K} \sum_{I=1}^K v_g^I$ .

3.  $v_g^I$  is the unadjusted sample variance within subtype  $I$ , i.e.  $v_g^I = \frac{1}{n^I} \sum_{s,z_s=I} (X_{g,s} - m_g^I)^2$  with  $m_g^I$  is the sample mean within subtype  $I$  and  $n^I$  is the number of cells within subtype  $I$ .

We estimate pooled over-dispersion rate by  $o_g = \frac{v_g}{m_g}$  and obtain  $\sigma_g = m_g \frac{o_g}{1-o_g}$  by the formula of mean of NB

What we are interested at those  $K$  groups comparison is the expression pattern,

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

For example  $K = 3$ , there are 5 expression pattern,  $P_1, P_2, \dots, P_5$ , Comparison between  $\mu$  is equivalent as comparison between  $q$ .

$$\begin{aligned} P1 : q_g^1 &= q_g^2 = q_g^3 \\ P2 : q_g^1 &= q_g^2 \neq q_g^3 \\ P3 : q_g^1 &\neq q_g^2 = q_g^3 \\ P4 : q_g^1 &= q_g^3 \neq q_g^2 \\ P5 : q_g^1 &\neq q_g^2 \neq q_g^3 \text{ and } q_g^1 \neq q_g^3 \end{aligned}$$

Under the assumption that two groups  $I$  and  $J$  share the same  $q_g$  we can pool the counts from the two groups by viewing them come from same distribution i.e.  $X_g^{I,J} | \sigma_g, q_g \sim NB(\sigma_g, q_g)$ ,  $q_g | \alpha, \beta^g \sim \text{Beta}(\alpha, \beta^g)$  and obtained the prior predictive function  $f(X_g^{I,J}) = \int_0^1 P(X_g^{I,J} | r_g, q_g) * P(q_g | \alpha, \beta^g) dq_g = \left[ \prod_{s=1}^S \binom{X_{g,s} + \sigma_g - 1}{X_{g,s}} \right] \frac{\text{Beta}(\alpha + \sum_{s=1}^S \sigma_g, \beta^g + \sum_{s=1}^S X_{g,s})}{\text{Beta}(\alpha, \beta^g)}$ . Consequently, we have prior predictive function for  $P1, \dots, P5$  as

$$\begin{aligned} h_1^g(X_g) &= f(X_g^{1,2,3}) \\ h_2^g(X_g) &= f(X_g^{1,2})f(X_g^3) \\ h_3^g(X_g) &= f(X_g^1)f(X_g^{2,3}) \\ h_4^g(X_g) &= f(X_g^{1,3})f(X_g^2) \\ h_5^g(X_g) &= f(X_g^1)f(X_g^2)f(X_g^3) \end{aligned}$$

where  $h_i^g(X_g) = P(X_g | M_{g,\pi_i}, z)$  for associated  $\pi_i$ . Then the marginal distribution of counts  $X_g$  is  $\sum_{k=1}^5 p_k h_k^g(X_g)$ , where the marginal  $p_k = P(M_{g,\pi} | z)$  (shared by all genome). Thus, the posterior probability of an expression pattern  $k$  is obtained by:

$$\frac{p_k h_k^g(X_g)}{\sum_{k=1}^5 p_k h_k^g(X_g)}$$

In the optimization steps for determining the hyper parameters  $(\alpha, \beta^g, p)$ , the computation and memory increase exponentially with the number of subtypes  $K$ . We use one-step EM as an approximation for the solution, that is  $\alpha$  and  $\beta^g$  are updated through gradient ascent.  $p$  is updated by the explicit form of the maximizer of the log likelihood.

2.2.2. *modalClust*. In this section, we review the procedure of modal-clustering and proof the compatibility of our Poisson-Gamma model.

### Product Partition Model

Let  $X = (X_1, X_2, \dots, X_n)$  be  $n$  one dimension observed data, given a partition for the data  $\pi = \{S_1, \dots, S_q\}$ , where  $S_i$  are disjoint subsets of  $\{1, 2, \dots, n\}$  and  $\bigcup_{i=1}^q S_i = \{1, 2, \dots, n\}$ . The likelihood for  $X$  satisfying such partition is

$$p(X|\pi) = \prod_{i=1}^q f(X_{S_i})$$

where  $X_{S_i}$  is the vector of observations corresponding to the items of component  $S_i$ . The component likelihood  $f(X_S)$  is defined for any non-empty component  $S$  and can take any form. The partition  $\pi$  is the only parameter we are interested at. Any other parameters that may have been involved in the model have been integrated over their prior.

The prior distribution for a partition  $\pi$  is also taken as a product form. We use the MAP partition (maximize the posterior  $p(\pi|X) \propto p(X|\pi)p(\pi)$ ) as the estimated clustering.

Dahl demonstrated by some choice of  $f$  and prior of  $\pi$ , we can reduce the time complexity of finding the MAP partition from factorial( $n$ ) to  $O(n^2)$  (Dahl, 2009). And the crucial condition for  $f$  is that if  $X_{S_1}$  and  $X_{S_2}$  are overlapped in the sense that  $\min\{X_{S_2}\} < \max\{X_{S_1}\} < \max\{X_{S_2}\}$  or  $\min\{X_{S_1}\} < \max\{X_{S_2}\} < \max\{X_{S_1}\}$ , let  $X_{S_1^*}$  and  $X_{S_2^*}$  be the sets of swapping one pair of those overlapped terms and keep the other unchanged. Then  $f(X_{S_1})f(X_{S_2}) \leq f(X_{S_1^*})f(X_{S_2^*})$ . Under such condition, we know that possible MAP candidates must be those partition  $\pi$ s that for any two blocks  $b_1, b_2 \in \pi$ , either  $\max_{i \in b_1}(X_i) \leq \min_{j \in b_2}(X_j)$  or  $\min_{i \in b_1}(X_i) \geq \max_{j \in b_2}(X_j)$ .

In Poisson-Gamma Model we assuming:

$$\begin{aligned} X_i|\pi, \lambda &\sim \text{Poisson}(X_i|\lambda_1 \mathbf{I}\{i \in S_1\} + \dots + \lambda_q \mathbf{I}\{i \in S_q\}) \\ \pi &\sim p(\pi) \\ \lambda_j &\sim \text{Gamma}(\alpha_0, \beta_0) \end{aligned}$$

where  $p(\pi) \propto \prod_{i=1}^q \eta_0 \Gamma(|S_i|)$ . Integrate out  $\lambda$ ,  $f(X_S)$  is obtained as:

$$f(X_S) = \frac{\beta^\alpha}{(|S| + \beta)^{\sum_{i \in S} X_i + \alpha}} \frac{\Gamma(\sum_{i \in S} X_i + \alpha)}{\Gamma(\alpha)} \frac{1}{\prod_{i \in S} X_i}$$

To apply modal-cluster on Poisson-Gamma model, we need to show  $f(X_S)$  still satisfying the condition mentioned before.

PROOF. if  $X_{S_1}$  and  $X_{S_2}$  are overlapped, without loss of generality, we assume  $\min\{X_{S_2}\} < \max\{X_{S_1}\} < \max\{X_{S_2}\}$ , and we swap  $\max\{X_{S_1}\}$  with  $\min\{X_{S_2}\}$  and keep the rest unchanged or we could also swap  $\max\{X_{S_1}\}$  with  $\max\{X_{S_2}\}$ . We denote the new set forming by swap of  $\max\{X_{S_1}\}$  with  $\min\{X_{S_2}\}$  as  $S_1^*$  and  $S_2^*$  and swap of  $\max\{X_{S_1}\}$  with  $\max\{X_{S_2}\}$  as  $S_1^{**}, S_2^{**}$  accordingly.

Then we need to show at least one of the following happens

- (1)  $f(X_{S_1^*})f(X_{S_2^*}) \geq f(X_{S_1})f(X_{S_2})$
- (2)  $f(X_{S_1^{**}})f(X_{S_2^{**}}) \geq f(X_{S_1})f(X_{S_2})$

Let  $a = \max\{X_{S_1}\}$ ,  $b = \min\{X_{S_2}\}$  and  $c = \max\{X_{S_2}\}$ .  $h_1 = \sum_{i \in S_1} X_i - a$  and  $h_2 = \sum_{i \in S_2} X_i - b$ ,  $n_1$  and  $n_2$  are the number of elements in  $S_1$  and  $S_2$ . Then

$$\begin{aligned}
 f(X_{S_1^*})f(X_{S_2^*}) &\geq f(X_{S_1})f(X_{S_2}) \\
 &\iff \\
 \frac{\Gamma(h_1 + a + \alpha)}{(n_1 + \beta)^{h_1 + a + \alpha}} \frac{\Gamma(h_2 + b + \alpha)}{(n_2 + \beta)^{h_2 + b + \alpha}} &\leq \frac{\Gamma(h_2 + a + \alpha)}{(n_2 + \beta)^{h_2 + a + \alpha}} \frac{\Gamma(h_1 + b + \alpha)}{(n_2 + \beta)^{h_1 + b + \alpha}} \\
 &\iff \\
 \frac{\Gamma(h_1 + a + \alpha)}{\Gamma(h_1 + b + \alpha)} \frac{\Gamma(h_2 + b + \alpha)}{\Gamma(h_2 + a + \alpha)} &\leq \left(\frac{n_1 + \beta}{n_2 + \beta}\right)^{a-b}
 \end{aligned}$$

Left hand side of above formula is  $\text{LHS}_1 = \frac{(h_1+b+\alpha)\dots(h_1+a-1+\alpha)}{(h_2+b+\alpha)\dots(h_2+a-1+\alpha)}$  by the property of Gamma function and  $X_i$  are integer.

Similarly,

$$\begin{aligned}
 f(X_{S_1^{**}})f(X_{S_2^{**}}) &\geq f(X_{S_1})f(X_{S_2}) \\
 &\iff \\
 \frac{\Gamma(h_2 + c + \alpha)}{\Gamma(h_2 + a + \alpha)} \frac{\Gamma(h_1 + a + \alpha)}{\Gamma(h_1 + c + \alpha)} &\leq \left(\frac{n_2 + \beta}{n_1 + \beta}\right)^{c-a}
 \end{aligned}$$

Left hand side of above formula is  $\text{LHS}_2 = \frac{(h_2+a+\alpha)\dots(h_2+c-1+\alpha)}{(h_1+a+\alpha)\dots(h_1+c-1+\alpha)}$

If  $h_1 \leq h_2$ , then  $\text{LHS}_1 \leq \left(\frac{h_1+a-1+\alpha}{h_2+a-1+\alpha}\right)^{a-b}$  and  $\text{LHS}_2 \leq \left(\frac{h_2+c-1+\alpha}{h_1+c-1+\alpha}\right)^{a-b}$

So if  $\frac{h_1+a-1+\alpha}{h_2+a-1+\alpha} \leq \frac{n_1+\beta}{n_2+\beta}$  then (12) holds, if  $\frac{h_2+c-1+\alpha}{h_1+c-1+\alpha} \leq \frac{n_1+\beta}{n_2+\beta}$  then (13) holds

We multiply those two inequalities, we found that  $\frac{h_1+a-1+\alpha}{h_2+a-1+\alpha} * \frac{h_2+c-1+\alpha}{h_1+c-1+\alpha} = \frac{h_1+a-1+\alpha}{h_1+c-1+\alpha} * \frac{h_2+c-1+\alpha}{h_2+a-1+\alpha} \leq 1$  as  $c > a$  and  $h_1 \leq h_2$  But  $\frac{n_1+\beta}{n_2+\beta} * \frac{n_1+\beta}{n_2+\beta} = 1$ . At least one equality holds, consequently at least one of (12) and (13) holds.

Similar proof for the case  $h_1 > h_2$ .

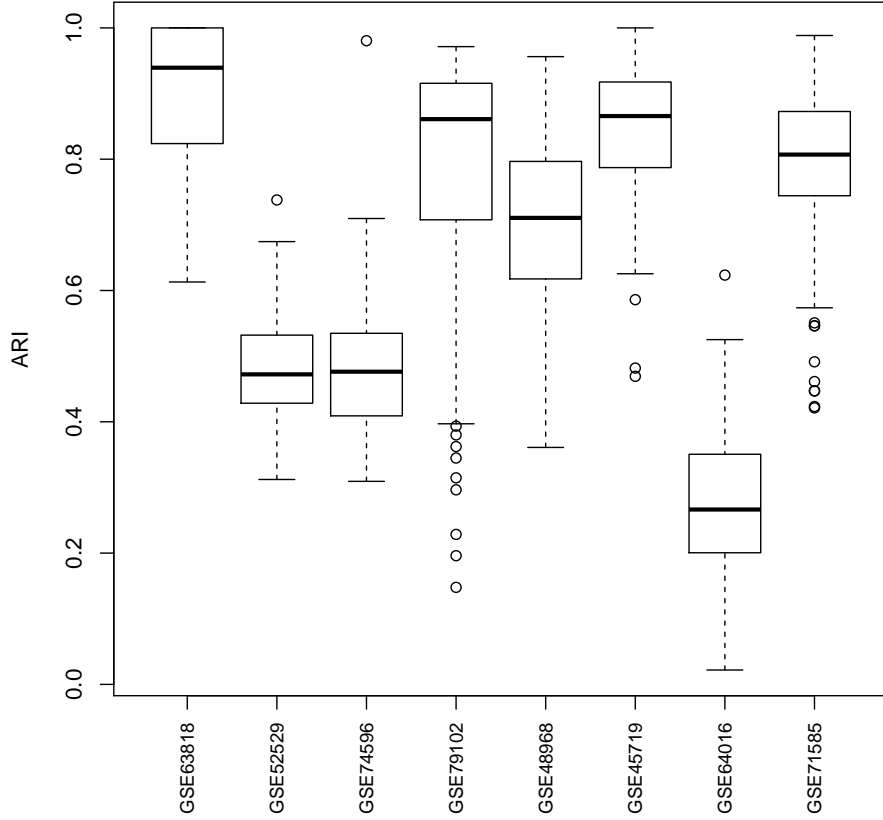
□

**2.2.3. Randomized K-means.** In this section, we illustrate how we estimated parameters for random weighing and demonstrate variability, accuracy and approximation to fully bayesian scheme of generated clustering.

To find the value of  $a_0, a_1$  and  $d_0$ , we have the marginal likelihood of  $d_{i,j}$ .

$$P(d_{i,j}|a_0, a_1, d_0) = \frac{\Gamma(a_0 + a_1)}{\Gamma(a_0)\Gamma(a_1)} \frac{d_0^{a_0} d_{i,j}^{a_1-1} a_1^{a_1}}{(d_0 + a_1 * d_{i,j})^{a_0+a_1}}$$

We estimate  $d_0$  by treating  $d_{i,j} \approx \Delta_{i,j}$  and based on the mean-variance ratio ( $\frac{E(1/\Delta_{i,j})}{\text{Var}(1/\Delta_{i,j})} = d_0$ ),  $d_0$  can be approximately estimated by moments of  $1/d_{i,j}$ , then we obtain  $a_0, a_1$  from MLE of marginal density of  $d_{i,j}$ . The MLE estimators are obtained through "nlminb" function in r, one issue is that the default



**Supplementary Figure S1:** Adjusted rand indexes to the clustering based on the original distance matrix without dividing weights. We investigate the randomness of clustering given by our weights through 8 datasets. All have stopping threshold for nlminb optimizing function in r with relative tolerance as 0.001

value for tolerance rate of stopping is  $1e-10$ , which yields large value of  $a_1 + a_0$  and resulting in non-randomness of our weighting matrix. We set tolerance rate as  $1e-3$ . And obtained moderate deviation from  $D$  (supplementary Figure S1)

We plot the ARI(adjusted random index) between the randomly generated clustering to clustering under the original distance across eight datasets. Though the mean varies, the interquartile range is wide enough presenting a reasonable variation of our random weighting scheme.

We also check validity of random weighting on simulated dataset. We random generate one-dimensional data  $X$  from a mixture of 5 normal distributions with different means and same variance. We compare clustering results between random weighting and bayesian clustering with Dirichlet process prior in terms of posterior probabilities that two elements belong to the same class given the whole data. We also compare accuracy of the two procedures by looking at the ARI comparing to true class label (supplementary Figure S2). We found that random weighting scheme tends to give better results than classical bayesian clustering.

**2.2.4. Selecting  $K$ .** In this section, we gave the criterion to select  $K$ .

In order to determine the number of clusters and inspired by validity defined in Ray and Turi (2000), We consider a modified validity =  $\frac{\text{intra}}{\text{inter}}$ . where  $\text{intra} = \frac{1}{N} \sum_{i=1}^K \sum_{x \in C_i} \|x - z_i\|^2$ ,  $\text{inter} = \text{mean}(\|z_i - z_j\|^2), i, j = 1, 2, \dots, K$  and  $z_i$  is the center (medoids) of cluster  $i$ . **intra** is the average of distance of a point to its corresponding cluster center, which measures the compactness of clusters. We made a



**Supplementary Figure S2:** comparison between random weighting scheme and bayesian clustering procedure in terms of posterior probabilities that two elements belong to the same class given the whole data and adjusted rand index comparing to the underlying true class label

small change here, in original paper **inter** was defined as minimum distance between medoids, we use average instead for the purpose of getting a smoother quantity. **inter** is the average distance of two cluster centers, which measures the separation between clusters. We want to have a small intra-cluster distance and a big inter-cluster distance, consequently we want to minimize the validity. From empirical study, we constantly observe a monotone decreasing relation between number of clusters and validity. However this quantity stabilize when  $K$  is sufficiently large. The stopping rule for searching  $K$  is when  $\text{validity}_K < \epsilon$  is satisfied. We set the default value of  $\epsilon$  to be 1. As we found DD analysis results to be most consistent with other scRNA method.

2.3. *Double Dirichlet Mixture.* In this section, we gave proofs for the properties and theorem for DDM in section 2.3 of main paper.

On the double Dirichlet masses, using notations in the main paper we have density functions:

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

with

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

and

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

Those computing units will serve as key components for proofing property 1 ~ 6 in section 2.3

Proof of property 1

PROOF. When  $\phi$  and  $\psi$  only satisfy the coarsest constraints:  $\sum_{i=1}^K \phi_i = \sum_{i=1}^K \psi_i = 1$ .  $\phi$  and  $\psi$  are independently Dirichlet distributed. When  $\phi$  and  $\psi$  satisfy finer constraints,  $P(\phi|\psi) \neq P(\phi)$  as there is some subsets  $b \neq \pi$  such that  $\sum_{i \in b} \phi = \sum_{i \in b} \psi$ . So  $\phi$  and  $\psi$  are dependent  $\square$

Proof of property 2

PROOF.  $E_\pi(\phi_k) = E_\pi(\phi_k | \Phi_b) E_\Phi(\Phi_b) = E_{\tilde{\phi}_b}(\tilde{\phi}_k) E_\Phi(\Phi_b)$  where  $b$  is the block containing subtype index  $k$ . As  $\tilde{\phi}_b \sim \text{Dirichlet}_{N(b)}[\alpha_b^1]$  and  $\Phi_\pi \sim \text{Dirichlet}_{N(\pi)}[\beta_\pi]$  We have  $E_{\tilde{\phi}_b}(\tilde{\phi}_k) = \frac{\alpha_k^1}{\sum_{k' \in b} \alpha_{k'}^1}$  and  $E_\Phi(\Phi_b) = \frac{\beta_b}{\sum_{b' \in \pi} \beta_{b'}}$ . Similarly we could proof the case for  $E_\pi(\psi_k)$   $\square$

Proof of property 3

PROOF.  $t^1/t_\pi^1$  is independent with  $t^2/t_\pi^2$  conditioning on  $t_\pi^1$  and  $t_\pi^2$  by the Neutrality property of Dirichlet distribution  $\square$

Proof of property 4

PROOF. For  $j = 1, 2$ , let  $T_b^j$  be the vector of  $t_k^j$  such that  $k \in b$ . Recall  $t_b^j = \sum_{k \in b} t_k^j$ . Without loss of generality, we consider the case condition  $j = 1$ .

At the support of  $p_\pi$ , for different blocks,  $T_b^1 | \tilde{\phi}_b$  are mutually independent. Then we have factorization:

$$p_\pi(t^1 | t_\pi^1, y) = \prod_{b \in \pi} (p(T_b^1 | t_b^1, y))$$

and right hand side prior predictive function can be obtained via integral out  $\tilde{\phi}_b$  given the prior Dirichlet $[\alpha_b^1]$  and  $p(T_b^1 | \tilde{\phi}_b)$  is multinomial( $\tilde{\phi}_b$ ) distributed.

$$\begin{aligned} p(T_b^1 | t_b^1, y) &= \int_{\tilde{\phi}_b} p(T_b^1 | \tilde{\phi}_b) p(\tilde{\phi}_b) d\tilde{\phi}_b \\ &= \left\{ \left[ \frac{\Gamma(t_b^1 + 1)}{\prod_{k \in b} \Gamma(t_k^1 + 1)} \right] \left[ \frac{\Gamma(\sum_{k \in b} \alpha_k^j)}{\prod_{k \in b} \Gamma(\alpha_k^j)} \right] \left[ \frac{\prod_{k \in b} \Gamma(\alpha_k^j + t_k^j)}{\Gamma(t_b^j + \sum_{k \in b} \alpha_k^j)} \right] \right\} \end{aligned}$$



□

### Proof of property 5

PROOF.  $t_\pi^1$  and  $t_\pi^2$  given the condition label  $y$  are independent identical distributed.  $t_\pi^1|\Phi \sim \text{multinomial}(\Phi)$

$$\begin{aligned} p_\pi(t_\pi^1, t_\pi^2|y) &= \int_{\Phi} p(t_\pi^1|\Phi) p(t_\pi^2|\Phi) p(\Phi) d\Phi \\ &= \left[ \frac{\Gamma(n_1+1)\Gamma(n_2+1)}{\prod_{b \in \pi} \Gamma(t_b^1+1)\Gamma(t_b^2+1)} \right] \left[ \frac{\Gamma(\sum_{b \in \pi} \beta_b)}{\prod_{b \in \pi} \Gamma(\beta_b)} \right] \left[ \frac{\prod_{b \in \pi} \Gamma(\beta_b + t_b^1 + t_b^2)}{\Gamma(n_1 + n_2 + \sum_{b \in \pi} \beta_b)} \right]. \end{aligned}$$

As prior of  $\Phi$  is Dirchlet $[\beta]$  and  $n_j = \sum_{b \in \pi} t_b^j$  for  $j = 1, 2$

□

To prove property 6, we need a lemma of dimensionality of the intersection of two  $A_{\pi}$ s

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 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  $N(\pi_2) + 1$ ,  $\dim(A_{\pi_1} \cap A_{\pi_2}) < \dim(A_{\pi_2})$

□

### Proof of property 6

PROOF. For a  $\pi$ ,  $P(A_\pi, |y, z) = \sum_{\tilde{\pi} \in \Pi} \int_{A_\pi} \omega_{\tilde{\pi}}^{\text{post}} d\phi d\psi$ , notice the support of  $\omega_{\tilde{\pi}}^{\text{post}}$  is  $A_{\tilde{\pi}}$ . By lemma 1, we know if  $\tilde{\pi}$  does not refine  $\pi$ , then  $\int_{A_\pi} \omega_{\tilde{\pi}}^{\text{post}} d\phi d\psi$  is an integral on lower dimension set and vanish. if  $\tilde{\pi}$  refines  $\pi$ , then  $\int_{A_\pi} \omega_{\tilde{\pi}}^{\text{post}} d\phi d\psi = \int_{A_{\tilde{\pi}}} \omega_{\tilde{\pi}}^{\text{post}} d\phi d\psi = \omega_{\tilde{\pi}}^{\text{post}}$ . We have  $P(A_\pi, |y, z) = \sum_{\tilde{\pi} \in \Pi} \omega_{\tilde{\pi}}^{\text{post}} 1[\tilde{\pi} \text{ refines } \pi]$

□

### Proof of theorem 3

PROOF. Recall the DDM prior:  $p(\phi, \psi) = \sum_{\pi \in \Pi} p_\pi(\phi, \psi)$ . By bayes' rule we know  $p(\phi, \psi|y, z) \propto p(\phi, \psi, y, z) = \sum_{\pi \in \Pi} p(y, z|\phi, \psi) p_\pi(\phi, \psi) \omega_\pi$  and we use the 1-1 map from  $(\phi, \psi)$  to  $(\tilde{\phi}, \tilde{\psi}, \Phi)$  to get

$$p(y, z|\phi, \psi) p_\pi(\phi, \psi) = p(y, z|\tilde{\phi}, \tilde{\psi}, \Phi_\pi) p(\tilde{\phi}) p(\tilde{\psi}) p(\Phi_\pi)$$

when  $(\phi, \psi) \in A_\pi$ . Let us denote right hand side of the above equation as  $U_\pi$ , then

$$U_\pi = \omega_\pi A_1 A_2 A_3 \prod_{k=1}^K (\tilde{\phi}_k)^{t_k^1 + \alpha_k^1} (\tilde{\psi}_k)^{t_k^2 + \alpha_k^2} \prod_{b \in \pi} (\Phi_b)^{t_b^1 + t_b^2 + \beta_b}$$

Where  $A_1$  is the product of normalizing terms from multinomial distribution of  $z^1$  and  $z^2$ ,  $A_1 = \frac{\Gamma(n_1+1)\Gamma(n_2+1)}{\prod_{j=1}^2 \prod_{k=1}^K \Gamma(t_k^j+1)}$

$A_2$  is the product of normalizing terms from Dirichlet distribution of  $\tilde{\phi}$  and  $\tilde{\psi}$ ,  $A_2 = \frac{\Gamma(\sum_{k=1}^K \alpha_k^1+1)\Gamma(\sum_{k=1}^K \alpha_k^2+1)}{\prod_{j=1}^2 \prod_{k=1}^K \Gamma(\alpha_k^j+1)}$

$A_3$  is the normalizing term from Dirichle distribution of  $\Phi_\pi$ ,  $A_3 = \frac{\Gamma(\sum_{b \in \pi} \beta_b+1)}{\prod_{b \in \pi} \Gamma(\beta_b+1)}$

Looking at the indexes of  $\tilde{\phi}, \tilde{\psi}$  and  $\Phi$ , we can decompose  $U_\pi$  as product of three Dirichlet densities with a normalizing term. Namely  $U_\pi = C_\pi * f_1 f_2 f_3$ , where  $f_1 \sim \text{Dirichlet}[\alpha^1 + t^1]$ ,  $f_2 \sim \text{Dirichlet}[\alpha^2 + t^2]$  and  $f_3 \sim \text{Dirichlet}[\beta + t^1 + t^2]$ . Considering the normalizing factors for densities  $f_1, f_2$  and  $f_3$ , and multiplying them with  $A_1, A_2$  and  $A_3$ . We have the  $C_\pi = p_\pi(t^1 | t_\pi^1, y) p_\pi(t^2 | t_\pi^2, y) p_\pi(t_\pi^1, t_\pi^2 | y) \omega_\pi$ . Consequently, we have

$$(\phi, \psi) | y, z \sim \text{DDM} \left[ \omega^{\text{post}} = (\omega_\pi^{\text{post}}), \alpha^1 + t^1, \alpha^2 + t^2 \right] \text{ and } \omega_\pi^{\text{post}} \propto p_\pi(t^1 | t_\pi^1, y) p_\pi(t^2 | t_\pi^2, y) p_\pi(t_\pi^1, t_\pi^2 | y) \omega_\pi.$$

Notice in DDM, we restricted  $\beta = \alpha^1 + \alpha^2$ .

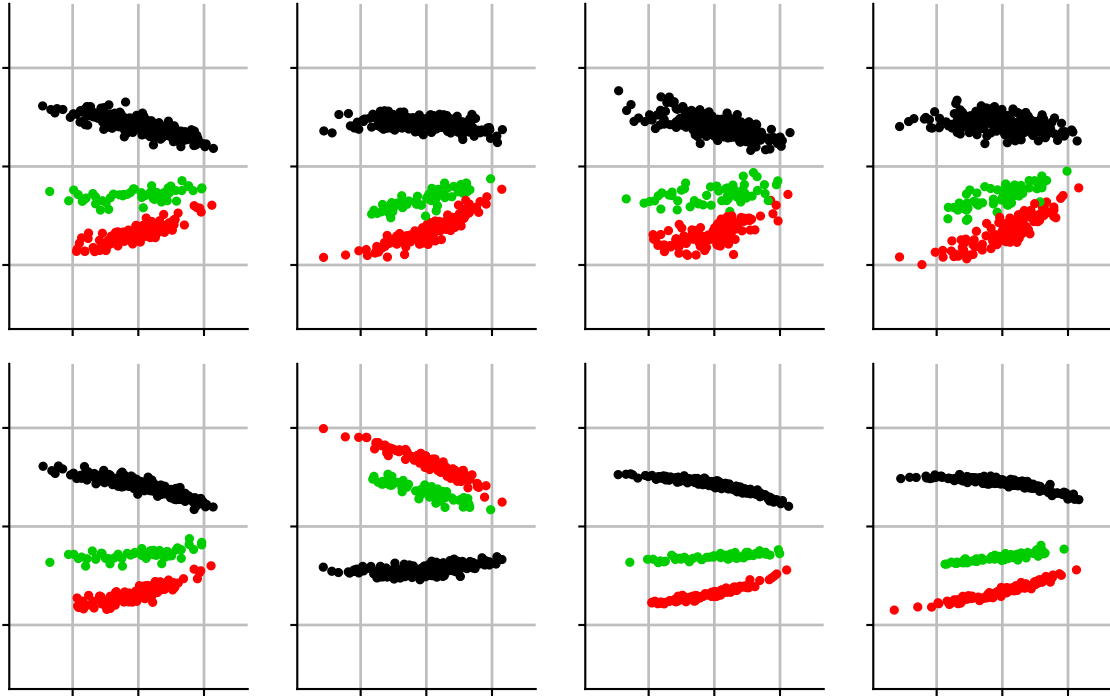
□

### 3. Numerical Experiments.

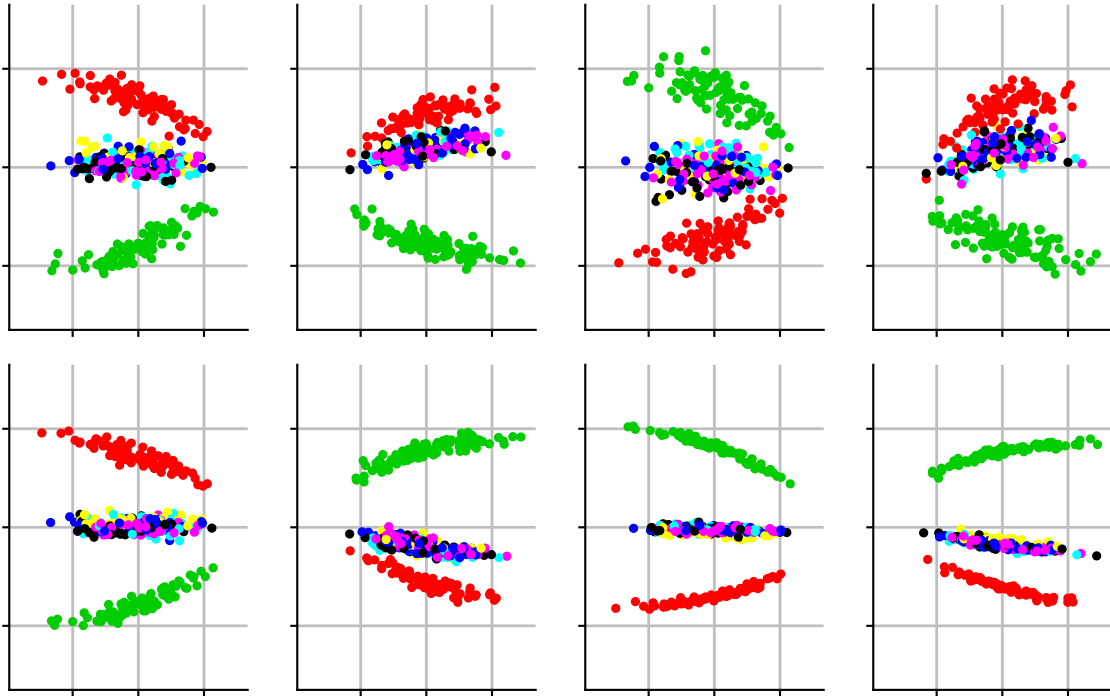
**3.1. Synthetic Data.** In this section, we use pca plots to show the subtle changes underlying each subtypes of simulated data and we demonstrate consistency of distributional changes based on scDDboost and Wasserstein distance between the empirical distribution of transcripts. Finally, roc curve illustrates that scDDboost having a good operating characteristic.

We first look at the pca plots of the simulated data (supplementary figure S3, S4, S5). For  $K = 7$  and 12, each scenario there were some subtypes nested in the 2d PCA projection. The distributional change of transcripts becomes difficult to detect. scDDboost benefits from the compositional structure and is more sensitive to those subtle changes.

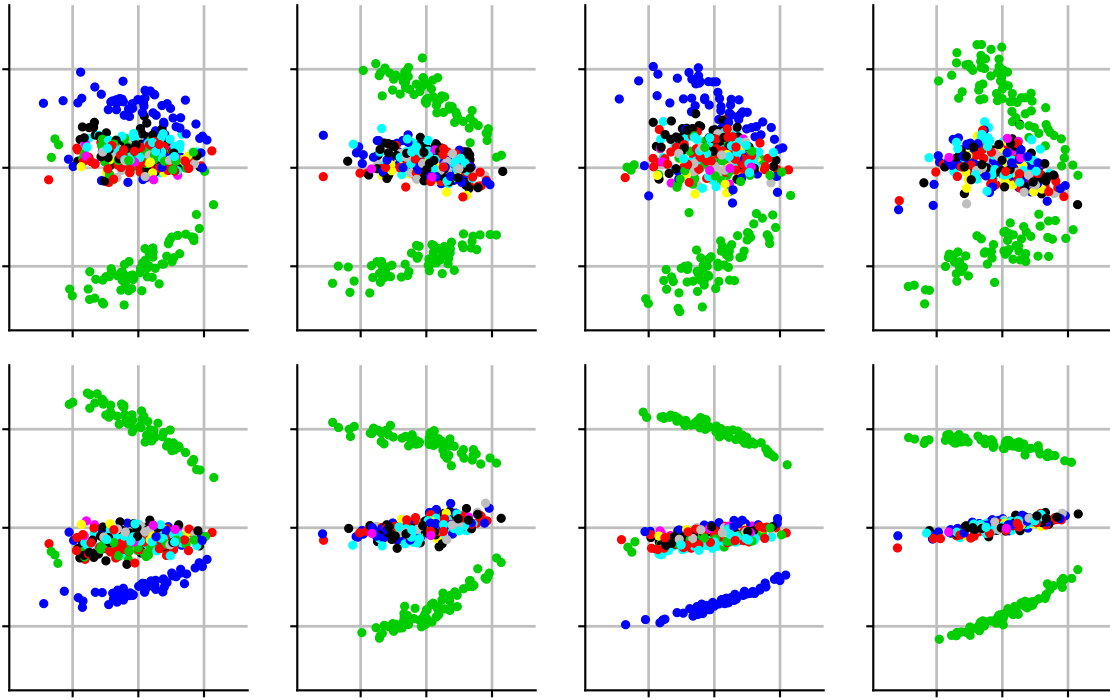
We observed consistent measurements of distributional change based on scDDboost and Wasserstein distance between the empirical distribution of transcripts.(supplementary figure S6) Lower probabilities of equivalent distributed are associated with bigger distances.



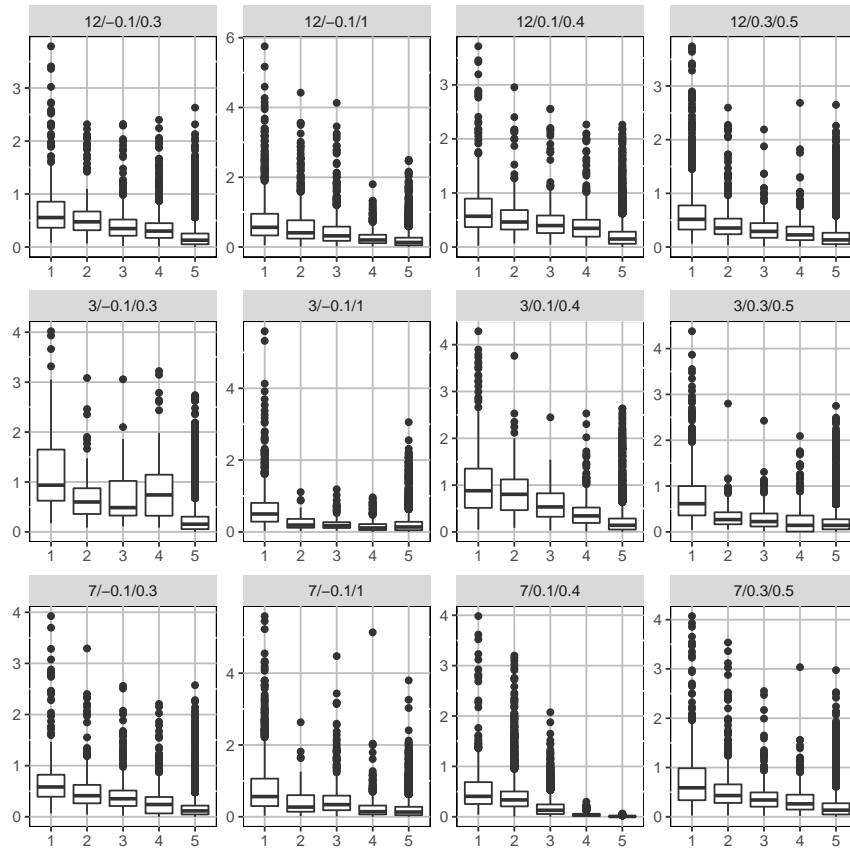
**Supplementary Figure S3:** first two principal components of transcripts under different parameters for simulated data. Different parameters resulted in different degree of separation of subtypes. We have 4 different settings for hyper-parameters of simulation, each setting has 2 replicates  $K = 3$



**Supplementary Figure S4:**  $K = 7$

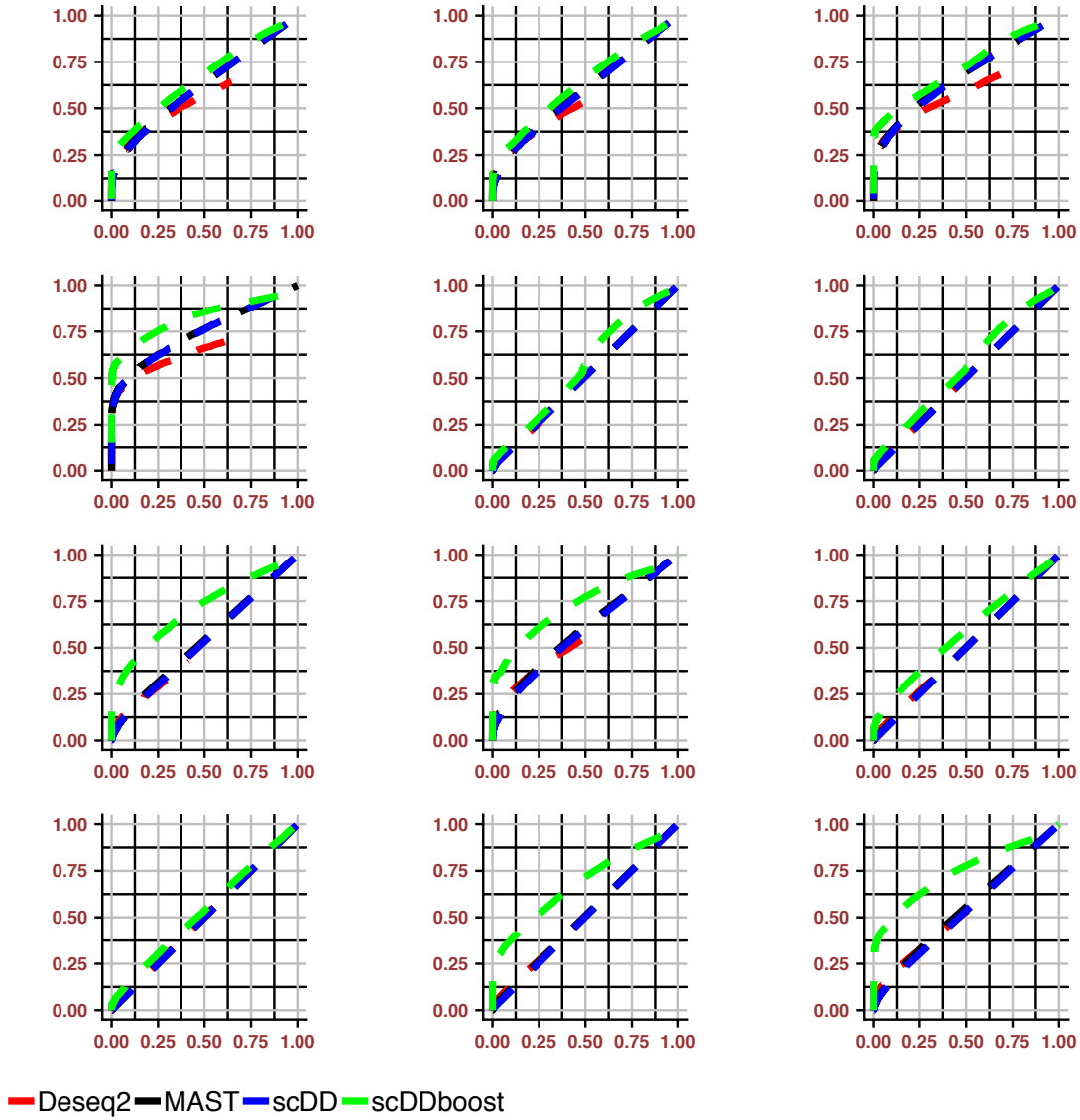


Supplementary Figure S5: K = 12



Supplementary Figure S6:  $P(ED_g|X, y)$  given by scDDboost versus empirical Wasserstein distance. Genes associated with boxes from left to right having  $P(ED_g|X, y)$  range from 0 - 0.2, 0.2 - 0.4, 0.4 - 0.6, 0.6 - 0.8, 0.8 - 1. For the simulation cases

We also have roc curve for the simulated data, each sub-figure is averaged over two replicates under the same parameters setting. scDDboost tends to outperform other methods (supplementary figure S7)



**Supplementary Figure S7:** Roc curve of the 12 simulation settings, under each setting, TPR and FPR are averaged over two replicates, generally we found scDDboost perform better than other methods

3.2. *Empirical Study.* In this section, we gave details of the empirical datasets and also demonstrate consistency to Wasserstein distance on one dataset (FUCCI).

**Data sets** details for the datasets used in the empirical studies of the main paper and the estimated number of subtypes  $K$  (supplementary table S1)

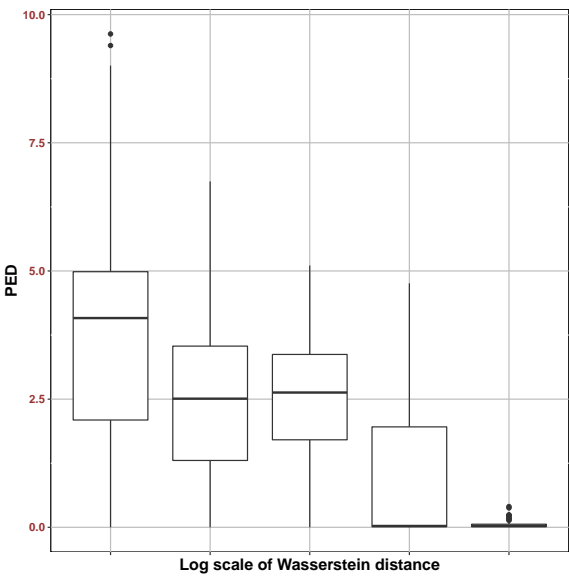
Data set	Conditions	Number of cells/condition	Organism	Ref	K
GSE94383	0 min unstim vs 75min stim	186,145	human	(Lane et al., 2017)	9
GSE48968-GPL13112	BMDC (2h LPS stimulation) vs 6h LPS	96,96	mouse	(Shalek et al., 2014)	4
GSE52529	T0 vs T72	69,74	human	(Trapnell et al., 2014)	7
GSE74596	NKT1 vs NTK2	46,68	mouse	(Engel et al., 2016)	7
EMTAB2805	G1 vs G2M	95,96	mouse	(Buettner et al., 2015)	6
GSE71585-GPL13112	Gad2tdTpositive vs Cux2tdTnegative	80,140	mouse	(Tasic et al., 2016)	4
GSE64016	G1 vs G2	91,76	human	(Leng et al., 2015)	6
GSE79102	patient1 vs patient2	51, 89	human	(Kiselev et al., 2017)	4
GSE45719	16-cell stage blastomere vs mid blastocyst cell	50, 60	mouse	(Deng et al., 2014)	4
GSE63818	Primordial Germ Cells, develop- mental stage: 7 week gestation vs Somatic Cells, developmental stage: 7 week gestation	40,26	mouse	(Guo et al., 2015)	6
GSE75748	DEC vs EC	64, 64	human	(Chu et al., 2016)	5
GSE84465	neoplastic cells vs non-neoplastic cells	1000, 1000	human	(Darmanis et al., 2017)	9

SUPPLEMENTARY TABLE S1  
datasets used for comparisons of DD analysis under different methods

The largest dataset we tried is GSE84465, we randomly sampled 1000 cells from each condition. The reason is that we in total have 3500 cells, it takes more time and memories to compute if we consider all the samples and 1000 cells each condition would be enough to represent the heterogeneity. We found DESeq identified significant smaller number of positives than others. It is intuitive that we are more likely to encounter subtle changes when we have large samples. Only consider mean shifts would have limited power. For other datasets we use all the cells within that condition under same batch.

We also observed consistent distributional change measurements by scDDboost and Wasserstein distance between the empirical distribution of transcripts . (supplementary Figure S8)

Datasets used for generating the Null cases (supplementary table S2)



**Supplementary Figure S8:**  $P(ED_g|X,y)$  given by scDDboost versus empirical Wasserstein distance. Genes associated with boxes from left to right having  $P(ED_g|X,y)$  range from 0 - 0.2, 0.2 - 0.4, 0.4 - 0.6, 0.6 - 0.8, 0.8 - 1, data used: FUCCI

Data set	Conditions	Number of cells/condition	Organism
GSE63818null	7 week gestation	20,20	mouse
GSE75748null	DEC	32, 32	human
GSE94383null	T0	93, 93	human
GSE48968-GPL13112null	BMDC (2h LPS stimulation)	48,48	mouse
GSE74596null	NKT1	23,23	mouse
EMTAB2805null	G1	48,48	mouse
GSE71585-GPL13112null	Gad2tdTpositive	40,40	mouse
GSE64016null	G1	46,45	human
GSE79102null	patient1	26, 25	human

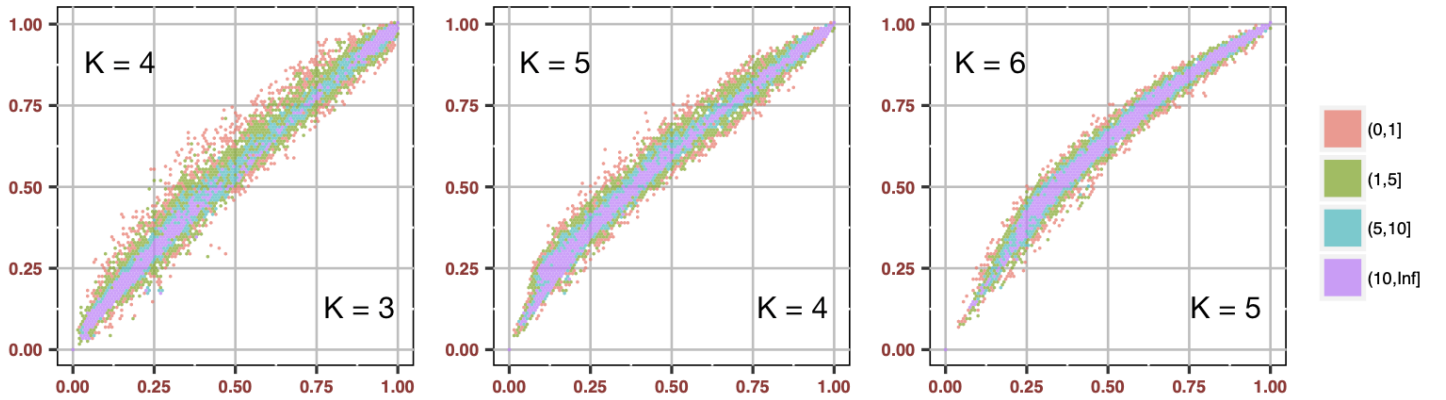
SUPPLEMENTARY TABLE S2

*datasets used for null cases, as cells are coming from same biological condition, there should not be any differential distributed genes, any positive call is false positive*

**3.3. Robustness.** In this section, we demonstrate change of PDD under different  $K$  and the robustness we gain through random weighting. We also gave a warning that using arbitrarily large  $K$  will inflate FDR.

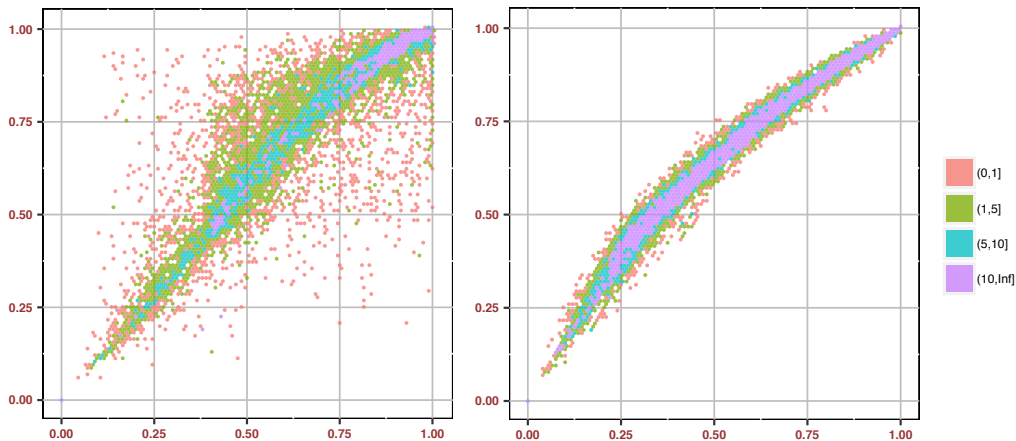
Number of subtypes  $K$  is a crucial factor controlling the accuracy of our modeling. Too small  $K$  may end up in an underfit such that cells within same subtype can still be very different, mean expression change among subtypes is incapable to capture the distribution change for some genes and consequently reducing the power of scDDboost. Too big  $K$  may end up in an overfit such that two subtypes can be very similar, given we have fixed number of samples (cells), allowing more clusters will introduce many patterns (both for mean expression change and proportion change) to infer. Also notice the limitation of DDM model (see section 4), overestimating  $K$  in scDDboost may lose FDR control (supplementary Figure S11).

From our empirical experience, it would be sufficient to capture the heterogeneity underlying cells with number of clusters less than 10 (supplementary Table 1). And we generally obtain stable validity score and PDD simultaneously (supplementary Figure S9). We demonstrate the change of PDD given different  $K$  at data GSE75748. When we increase  $K$ , the variance of the differential term  $PDD_{K+1} - PDD_K$  keeps decreasing and PDD keeps increasing. Our selection criterion ( $K = 5$ ) happens to choose  $K$  such that change between  $PDD_{K+1}$  and  $PDD_K$  is small while not inflating PDD.



**Supplementary Figure S9:** PDD change under different number of subtypes  $K$ , dataset used DEC-EC, our rule for selecting  $K$  tends also to make PDD stabilize

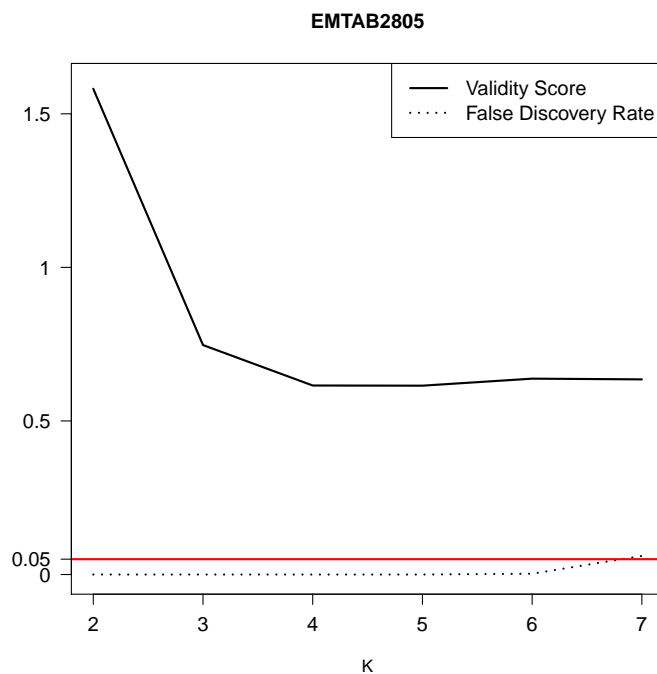
The random weighting scheme help us obtain robust PDD.



**Supplementary Figure S10:** DEC-EC, PDD under  $K = 5$  vs.  $K = 6$ , left panel is without the randomized distance and right panel is with randomized distance. We increase robustness of our methods through random weighting



scDDboost may lose FDR control if we keep increasing  $K$ .



**Supplementary Figure S11:** under NULL case, using dataset EMTAB2805, when using too big  $K$  we may lose FDR control (black dashed line shows proportion of false positive identified by scDDboost under 0.05 threshold, while validity score did not vary too much after  $K$  is greater than 2)

#### 4. Posterior consistency. In this section, we proof theorem 4.

As the density of DDM is computed by product or ratio over bunches of gamma function and gamma function is not easy to direct work on it and derive limiting theorem. To proof theorem 4, we need a crucial lemma which gave us an approximation to the gamma function, namely

LEMMA 2. For  $x \geq 1$ ,  $\frac{x^{x-c}}{e^{x-1}} \leq \Gamma(x) \leq \frac{x^{x-1/2}}{e^{x-1}}$ , where  $c = 0.577215\dots$  is the Euler-Mascheroni constant.

PROOF. By (Li and ping Chen, 2007), we have  $\frac{x^{x-c}}{e^{x-1}} \leq \Gamma(x) \leq \frac{x^{x-1/2}}{e^{x-1}}$  for  $x > 1$  and now we added the case when  $x = 1, \Gamma(x) = 1$  so that both sides will include the equality case.  $\square$

LEMMA 3. For positive integer  $n$ ,  $\sqrt{2\pi n^{n+1/2}}e^{-n} \leq \Gamma(n+1) \leq en^{n+1/2}e^{-n}$

We have another two lemmas and theorem 1 and 2 are just proporsition of the lemma

LEMMA 4. If  $(\phi, \psi) \in A_{\pi_1} \cap A_{\pi_2}$ , follow the conditions in theorem 1 then

$$\frac{\omega_{\pi_1}^{post}}{\omega_{\pi_2}^{post}} \xrightarrow[n \rightarrow \infty]{a.s.} 0 \quad \text{if } N(\pi_1) < N(\pi_2)$$

PROOF. Recall  $\omega_{\pi}^{post} \propto p_{\pi}(t^1|t_{\pi}^1, y) p_{\pi}(t^2|t_{\pi}^2, y) p_{\pi}(t_{\pi}^1, t_{\pi}^2|y) \omega_{\pi}$ . and  $\text{RHS} = g(\pi, \alpha, \beta, n_1, n_2) f(\pi, t^1, t^2, \alpha, \beta)$  and  $\frac{\omega_{\pi_1}^{post}}{\omega_{\pi_2}^{post}} = \frac{g(\pi_1, \alpha, \beta, n_1, n_2)}{g(\pi_2, \alpha, \beta, n_1, n_2)} \frac{f(\pi_1, t^1, t^2, \alpha, \beta)}{f(\pi_2, t^1, t^2, \alpha, \beta)}$  where

$$g(\pi, t^1, t^2, \alpha, \beta) = \left[ \prod_{j=1}^2 \prod_{b \in \pi} \frac{\Gamma(\sum_{k \in b} \alpha_k^j)}{\prod_{k \in b} \Gamma(\alpha_k^j)} \right] \frac{\Gamma(n_1 + 1) \Gamma(n_2 + 1)}{\prod_{b \in \pi} \Gamma(\beta_b)} \frac{\Gamma(\sum_{b \in \pi} \beta_b)}{\Gamma(n_1 + n_2 + \sum_{b \in \pi} \beta_b)}$$

$$f(\pi, t^1, t^2, \alpha, \beta) = \left[ \prod_{j=1}^2 \prod_{b \in \pi} \frac{1}{\prod_{k \in b} \Gamma(t_k^j + 1)} \frac{\prod_{k \in b} \Gamma(\alpha_k^j + t_k^j)}{\Gamma(t_b^j + \sum_{k \in b} \alpha_k^j)} \right] \prod_{b \in \pi} \Gamma(\beta_b + t_b^1 + t_b^2)$$

For notation simplicity, we use the abbreviation  $g(\pi), f(\pi)$  to substitute  $g(\pi, \alpha, \beta, n_1, n_2), f(\pi, t^1, t^2, \alpha, \beta)$ .

We take log on  $\frac{\omega_{\pi_1}^{post}}{\omega_{\pi_2}^{post}}$ , denote it as LR.  $\text{LR} = \ln g(\pi_1) - \ln g(\pi_2) + \ln f(\pi_1) - \ln f(\pi_2)$ . Denote  $C(\pi_1, \pi_2, \alpha, \beta) = \ln g(\pi_1) - \ln g(\pi_2)$ ,  $C(\pi_1, \pi_2, \alpha, \beta)$  does not change with sample size  $n_1, n_2$  and is a constant determined by partition  $\pi_1, \pi_2$  and hyper parameters  $\alpha, \beta$ . For further convenience of notation let  $h(x) = \ln \Gamma(x)$  and  $\gamma_b^j = \sum_{k \in b} \alpha_k^j$ . Denote  $R(\pi_1, \pi_2, t^1, t^2, \alpha, \beta) = \ln f(\pi_1) - \ln f(\pi_2)$ . And removing the common part of  $f(\pi_1)$  and  $f(\pi_2)$ , we have

$$R(\pi_1, \pi_2, t^1, t^2, \alpha, \beta) = d(\pi_1, t^1, t^2, \alpha, \beta) - d(\pi_2, t^1, t^2, \alpha, \beta)$$

where

$$d(\pi, t^1, t^2, \alpha, \beta) = \sum_{b \in \pi} h(\beta_b + t_b^1 + t_b^2) - \sum_{j=1}^2 \sum_{b \in \pi} h(t_b^j + \gamma_b^j)$$

Recall  $\beta_b = \gamma_b^1 + \gamma_b^2$  and from lemma 2,  $(x - c)\ln(x) - x + 1 \leq h(x) \leq (x - 1/2)\ln(x) - x + 1$  we have

$$(3) \quad d(\pi, t^1, t^2, \alpha, \beta) \geq \sum_{b \in \pi} (\beta_b + t_b^1 + t_b^2 - c) \ln(\beta_b + t_b^1 + t_b^2) - \sum_{j=1}^2 \sum_{b \in \pi} (t_b^j + \gamma_b^j - 1/2) \ln(t_b^j + \gamma_b^j) + N(\pi)$$

$$(4) \quad d(\pi, t^1, t^2, \alpha, \beta) \leq \sum_{b \in \pi} (\beta_b + t_b^1 + t_b^2 - 1/2) \ln(\beta_b + t_b^1 + t_b^2) - \sum_{j=1}^2 \sum_{b \in \pi} (t_b^j + \gamma_b^j - c) \ln(t_b^j + \gamma_b^j) + N(\pi)$$

$$\begin{aligned} \text{RHS of (4)} &= \Sigma_b [(t_b^1 + \gamma_b^1) \ln(1 + \frac{t_b^2 + \gamma_b^2}{t_b^1 + \gamma_b^1}) + (t_b^2 + \gamma_b^2) \ln(1 + \frac{t_b^1 + \gamma_b^1}{t_b^2 + \gamma_b^2}) \\ &+ (1 - c) \ln(\beta_b + t_b^1 + t_b^2) - 1/2(\ln(1 + \frac{t_b^2 + \gamma_b^2}{t_b^1 + \gamma_b^1}) + \ln(1 + \frac{t_b^1 + \gamma_b^1}{t_b^2 + \gamma_b^2}))] + N(\pi) \end{aligned}$$

By Taylor expansion at  $x = 1$ ,  $\ln(x + 1) = \ln 2 + 1/2(x - 1) - 1/8(x - 1)^2 + g(\xi)(x - 1)^3$ , where  $g(\xi)$  is the reminder term of form  $\frac{1}{3(1+\xi)^3}$  for  $0 < \xi < x$ . For a fixed  $n_1, n_2$ , we have

$$\begin{aligned} \text{RHS of (4)} &= (n_1 + n_2) \ln 2 - \Sigma_{b \in \pi} (1/8(X_b^1 + X_b^2) \\ &+ g(\xi_b)(Y_b^1 + Y_b^2)) + T(\pi) + N(\pi) \end{aligned}$$

where  $X_b^1 = \frac{(t_b^1 - t_b^2 + \gamma_b^1 - \gamma_b^2)^2}{t_b^1 + \gamma_b^1}$ ,  $X_b^2 = \frac{(t_b^1 - t_b^2 + \gamma_b^1 - \gamma_b^2)^2}{t_b^2 + \gamma_b^2}$ ,  $Y_b^1 = \frac{(t_b^1 - t_b^2 + \gamma_b^1 - \gamma_b^2)^3}{(t_b^1 + \gamma_b^1)^2}$ ,  $Y_b^2 = \frac{(t_b^1 - t_b^2 + \gamma_b^1 - \gamma_b^2)^3}{(t_b^2 + \gamma_b^2)^2}$  and  $T(\pi) = \Sigma_{b \in \pi} [(1 - c) \ln(\beta_b + t_b^1 + t_b^2) - 1/2(\ln(1 + \frac{t_b^2 + \gamma_b^2}{t_b^1 + \gamma_b^1}) + \ln(1 + \frac{t_b^1 + \gamma_b^1}{t_b^2 + \gamma_b^2}))]$

Similarly

$$\begin{aligned} \text{RHS of (5)} &= (n_1 + n_2) \ln 2 - \Sigma_{b \in \pi} (1/8(X_b^1 + X_b^2) \\ &+ g(\xi_b)(Y_b^1 + Y_b^2)) + U(\pi) + N(\pi) \end{aligned}$$

$$U(\pi) = \Sigma_{b \in \pi} [(2c - 1/2) \ln(\beta_b + t_b^1 + t_b^2) - c(\ln(1 + \frac{t_b^2 + \gamma_b^2}{t_b^1 + \gamma_b^1}) + \ln(1 + \frac{t_b^1 + \gamma_b^1}{t_b^2 + \gamma_b^2}))]$$

Using above inequalities, we have

$$\begin{aligned} R(\pi_1, \pi_2, t^1, t^2, \alpha, \beta) &\leq U(\pi_1) - T(\pi_2) - 1/8(\Sigma_{b \in \pi_1} (X_b^1 + X_b^2) - \Sigma_{b \in \pi_2} (X_b^1 + X_b^2)) \\ &+ \Sigma_{b \in \pi_1} g(\xi_b)(Y_b^1 + Y_b^2) - \Sigma_{b \in \pi_2} g(\xi_b)(Y_b^1 + Y_b^2) \end{aligned}$$

$Y_b^j = \frac{((t_b^1 - t_b^2 + \gamma_b^1 - \gamma_b^2)/\sqrt{n})^3/\sqrt{n}}{((t_b^1 + \gamma_b^1)/n)^2}$ , by LLN the denominator goes to a constant and by CLT in the numerator  $(t_b^1 - t_b^2 + \gamma_b^1 - \gamma_b^2)/\sqrt{n} \rightarrow (t_b^1 - t_b^2)/\sqrt{n} \rightarrow \sqrt{n}[(t_b^1/n - \Phi_b) - (t_b^2/n - \Psi_b)]$ , which goes to a normal distributed random variables when  $\Phi_b = \Psi_b$ . So  $Y_b^j$  is  $o_p(1)$ . Similarly,  $X_b^j = \frac{((t_b^1 - t_b^2 + \gamma_b^1 - \gamma_b^2)/\sqrt{n})^2}{t_b^j + \gamma_b^j/n}$  is asymptotic gamma( $\chi$ -square) distributed.  $g(\xi_b)$  has bounded variance,  $U(\pi_1) - T(\pi_2) = -\ln(n)$  if  $N(\pi_2) < N(\pi_1)$  as  $\ln(\beta_b + t_b^1 + t_b^2) - \ln(\beta_{b'} + t_{b'}^1 + t_{b'}^2) = \ln(\frac{\beta_b + t_b^1 + t_b^2}{n}) - \ln(\frac{\beta_{b'} + t_{b'}^1 + t_{b'}^2}{n}) \rightarrow O(1)$  a.s. so we complete the proof □

**LEMMA 5.** If  $(\phi, \psi) \in A_{\pi_1} \cap A_{\pi_2}$ , follow the conditions in theorem 1 and further we have  $\omega^j, j = 1, 2$  be vectors of integers then

$$\frac{\omega_{\pi_1}^{post}}{\omega_{\pi_2}^{post}} \xrightarrow[n \rightarrow \infty]{d} v \quad \text{if } N(\pi_1) = N(\pi_2)$$

$v$  is a random variable

PROOF. follow almost same procedure in lemma 4, but instead of using inequalities in lemma 2, we use lemma 3. And we still have

$$d(\pi, t^1, t^2, \alpha, \beta) = \sum_{b \in \pi} h(\beta_b + t_b^1 + t_b^2) - \sum_{j=1}^2 \sum_{b \in \pi} h(t_b + \gamma_b^j)$$

and by lemma 3

(5)

$$d(\pi, t^1, t^2, \alpha, \beta) \geq \sum_{b \in \pi} (\beta_b + t_b^1 + t_b^2 - 1/2) \ln(\beta_b + t_b^1 + t_b^2) - \sum_{j=1}^2 \sum_{b \in \pi} (t_b^j + \gamma_b^j - 1/2) \ln(t_b^j + \gamma_b^j) + \ln(\sqrt{2\pi}) - 1$$

(6)

$$d(\pi, t^1, t^2, \alpha, \beta) \leq \sum_{b \in \pi} (\beta_b + t_b^1 + t_b^2 - 1/2) \ln(\beta_b + t_b^1 + t_b^2) - \sum_{j=1}^2 \sum_{b \in \pi} (t_b^j + \gamma_b^j - 1/2) \ln(t_b^j + \gamma_b^j) + 1 - \ln(\sqrt{2\pi})$$

$$R(\pi_1, \pi_2, t^1, t^2, \alpha, \beta) \approx D(\pi_1) - D(\pi_2) - 1/8(\sum_{b \in \pi_1} (X_b^1 + X_b^2) - \sum_{b \in \pi_2} (X_b^1 + X_b^2)) \\ - \sum_{b \in \pi_1} g(\xi_b)(Y_b^1 + Y_b^2) - \sum_{b \in \pi_2} g(\xi_b)(Y_b^1 + Y_b^2)$$

where  $D(\pi) = \sum_{b \in \pi} [1/2 \ln(\beta_b + t_b^1 + t_b^2) - c(\ln(1 + \frac{t_b^2 + \gamma_b^2}{t_b^1 + \gamma_b^1}) + \ln(1 + \frac{t_b^1 + \gamma_b^1}{t_b^2 + \gamma_b^2}))]$  And  $D(\pi_1) - D(\pi_2)$  is  $O(1)$  if  $N(\pi_1) = N(\pi_2)$  as  $\ln(\beta_b + t_b^1 + t_b^2) - \ln(\beta_{b'} + t_{b'}^1 + t_{b'}^2) = \ln(\frac{\beta_b + t_b^1 + t_b^2}{n_1}) - \ln(\frac{\beta_{b'} + t_{b'}^1 + t_{b'}^2}{n_1}) \rightarrow 0 \quad a.s.$   $\square$

## Proof of theorem 4

PROOF. Recall  $\sum_{\pi \in \Pi} \omega_{\pi}^{\text{post}} = 1$  and  $P(A_{\pi}|y, z) = \sum_{\tilde{\pi} \in \Pi} \omega_{\tilde{\pi}}^{\text{post}} 1[\tilde{\pi} \text{ refines } \pi]$ . If  $(\phi, \psi) \notin Q$ , for all the  $A_{\pi}$  covers  $(\phi, \psi)$  there is one finest  $\pi^*$  with the largest  $N(\pi^*)$  and every other  $\pi$  that  $(\phi, \psi) \in A_{\pi}$  is coarser than  $\pi^*$ . We get the results of theorem 4 by lemma 4.  $\square$

1

## References.

- BUETTNER, F., NATARAJAN, K. N., CASALE, F. P., PROSERPIO, V., SCIALDONE, A., THEIS, F. J., TEICHMANN, S. A., MARIONI, J. C. and STEGLE, O. (2015). Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells. *Nature Biotechnology* **33** 155 EP -.
- CHU, L.-F., LENG, N., ZHANG, J., HOU, Z., MAMOTT, D., VEREIDE, D. T., CHOI, J., KENDZIORSKI, C., STEWART, R. and THOMSON, J. A. (2016). Single-cell RNA-seq reveals novel regulators of human embryonic stem cell differentiation to definitive endoderm. *Genome Biology* **17** 173. .
- DAHL, D. B. (2009). Modal clustering in a class of product partition models. *Bayesian Anal.* **4** 243–264.
- DARMANIS, S., SLOAN, S. A., CROOTE, D., MIGNARDI, M., CHERNIKOVA, S., SAMGHABABI, P., ZHANG, Y., NEFF, N., KOWARSKY, M., CANEDA, C., LI, G., CHANG, S. D., CONNOLLY, I. D., LI, Y., BARRES, B. A., GEPHART, M. H. and QUAKE, S. R. (2017). Single-Cell RNA-Seq Analysis of Infiltrating Neoplastic Cells at the Migrating Front of Human Glioblastoma. *Cell reports* **21** 1399–1410.
- DENG, Q., RAMSKÖLD, D., REINIUS, B. and SANDBERG, R. (2014). Single-Cell RNA-Seq Reveals Dynamic, Random Monoallelic Gene Expression in Mammalian Cells. *Science* **343** 193–196.
- ENGEL, I., SEUMOIS, G., CHAVEZ, L., SAMANIEGO-CASTRUITA, D., WHITE, B., CHAWLA, A., MOCK, D., VIJAYANAND, P. and KRONENBERG, M. (2016). Innate-like functions of natural killer T cell subsets result from highly divergent gene programs. *Nature Immunology* **17** 728 EP -.
- GUO, F., YAN, L., GUO, H., LI, L., HU, B., ZHAO, Y., YONG, J., HU, Y., WANG, X., WEI, Y., WANG, W., LI, R., YAN, J., ZHI, X., ZHANG, Y., JIN, H., ZHANG, W., HOU, Y., ZHU, P., LI, J., ZHANG, L., LIU, S., REN, Y., ZHU, X., WEN, L., GAO, Y. Q., TANG, F. and QIAO, J. (2015). The Transcriptome and DNA Methylome Landscapes of Human Primordial Germ Cells. *Cell* **161** 1437–1452.

- KISELEV, V. Y., KIRSCHNER, K., SCHAUB, M. T., ANDREWS, T., YIU, A., CHANDRA, T., NATARAJAN, K. N., REIK, W., BARAHONA, M., GREEN, A. R. and HEMBERG, M. (2017). SC3: consensus clustering of single-cell RNA-seq data. *Nature Methods* **14** 483 EP -.
- LANE, K., VAN VALEN, D., DEFELICE, M. M., MACKLIN, D. N., KUDO, T., JAIMOVICH, A., CARR, A., MEYER, T., PE'ER, D., BOUTET, S. C. and COVERT, M. W. (2017). Measuring Signaling and RNA-Seq in the Same Cell Links Gene Expression to Dynamic Patterns of NF-B Activation. *Cell Systems* **4** 458–469.e5.
- LENG, N., CHU, L.-F., BARRY, C., LI, Y., CHOI, J., LI, X., JIANG, P., STEWART, R. M., THOMSON, J. A. and KENDZIORSKI, C. (2015). Oscop identifies oscillatory genes in unsynchronized single-cell RNA-seq experiments. *Nature Methods* **12** 947 EP -.
- LI, X. and PING CHEN, C. (2007). Inequalities for the gamma function. In 2007), Art. 28. [ONLINE: <http://jipam.vu.edu.au/article.php?sid=842>].
- RAY, S. and TURI, R. H. (2000). Determination of Number of Clusters in K-Means Clustering and Application in Colour Image Segmentation.
- SHALEK, A. K., SATIJA, R., SHUGA, J., TROMBETTA, J. J., GENNERT, D., LU, D., CHEN, P., GERTNER, R. S., GAUBLomme, J. T., YOSEF, N., SCHWARTZ, S., FOWLER, B., WEAVER, S., WANG, J., WANG, X., DING, R., RAYCHOWDHURY, R., FRIEDMAN, N., HACOEN, N., PARK, H., MAY, A. P. and REGEV, A. (2014). Single-cell RNA-seq reveals dynamic paracrine control of cellular variation. *Nature* **510** 363 EP -.
- TASIC, B., MENON, V., NGUYEN, T. N., KIM, T. K., JARSKY, T., YAO, Z., LEVI, B., GRAY, L. T., SORESENSEN, S. A., DOLBEARE, T., BERTAGNOLLI, D., GOLDY, J., SHAPOVALOVA, N., PARRY, S., LEE, C., SMITH, K., BERNARD, A., MADISEN, L., SUNKIN, S. M., HAWRYLYCZ, M., KOCH, C. and ZENG, H. (2016). Adult mouse cortical cell taxonomy revealed by single cell transcriptomics. *Nature Neuroscience* **19** 335 EP -.
- TRAPNELL, C., CACCHIARELLI, D., GRIMSBY, J., POKHAREL, P., LI, S., MORSE, M., LENNON, N. J., LIVAK, K. J., MIKKELSEN, T. S. and RINN, J. L. (2014). The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nature biotechnology* **32** 381–386.