

# Supplementary note for “A deconvolution framework that uses single-cell sequencing plus a small benchmark dataset for accurate analysis of cell type ratios in complex tissue samples”

This note provides further mathematical details of the DeMixSC deconvolution framework.

## 1 Notation

Throughout, we use boldfaced lowercase for vectors, e.g.,  $\mathbf{s}$  denotes a vector of cell size for each cell type, and uppercase letters for matrices, e.g.,  $\mathbf{R}$  denotes the cell type specific reference matrix. We use plain lowercase and uppercase letters to denote elements and sets, respectively. For instance,  $g \in G$  means gene  $g$  belongs to a considered set of genes  $G$ . To distinguish the observed expression gene values from the underlying ground-truth objects, we add tildes to the latter. For example,  $y$  and  $x$  are the observed expression values of bulk and sc/snRNA-seq data, respectively;  $\tilde{y}$  and  $\tilde{x}$  are the true mean expressions. The estimated parameters are marked with “hat”, e.g.,  $\hat{p}$  denotes the estimated cell type proportion.

## 2 The DeMixSC framework

### 2.1 The wNNLS model

DeMixSC is a reference-based deconvolution framework built upon the widely adopted weighted non-negative least squares (wNNLS) model [1, 2, 3, 4] with multiple critical enhancements. To begin with, let  $j$  be the subject index,  $g$  be the gene index,  $c$  be the cell index, and  $k$  be the cell type index. We use  $G$  for denoting the set of genes,  $K$  for the set of cell types, and  $C_j^k$  for the set of cells of cell type  $k$  in subject  $j$ . The true expression of gene  $g$  in subject  $j$  at the bulk tissue level,  $\tilde{y}_{jg}$ , can be written as the summation of its true expressions at the single-cell level,  $\tilde{x}_{jgc}^k$ , as follows:

$$\tilde{y}_{jg} = \sum_{k \in K} \sum_{c \in C_j^k} \tilde{x}_{jgc}^k.$$

Note here that in this work, we use relative abundance as the gene expression of the bulk data as in MuSiC [4]. To account for measurement errors, we assume the observed gene expression at the bulk level  $y_{jg}$  as

$$y_{jg} = \tilde{y}_{jg} + \epsilon_{jg} = \sum_{k \in K} \sum_{c \in C_j^k} \tilde{x}_{jgc}^k + \epsilon_{jg}, \quad (\text{S.1})$$

where  $\epsilon_{jg}$  is the measurement noise for gene  $g$  in subject  $j$  introduced by the bulk sequencing platform, with mean zero and variance  $\sigma_{jg}^2$ .

Next, let  $n_j^k$  be the number of cells of cell type  $k$  in subject  $j$  and the number of all cells  $n_j = \sum_{k \in K} n_j^k$ . The cell type proportion of cell type  $k$  in subject  $j$  is then  $p_j^k = \frac{n_j^k}{n_j}$ . Then, the observed bulk gene

expression (S.1) can be expressed as

$$y_{jg} = n_j \sum_{k \in K} p_j^k \theta_{jg}^k s_j^k + \epsilon_{jg}, \quad (\text{S.2})$$

where  $\theta_{jg}^k = \frac{\sum_{c \in C_j^k} \tilde{x}_{jgc}^k}{\sum_{g \in G} \sum_{c \in C_j^k} \tilde{x}_{jgc}^k}$  is the relative abundance of gene  $g$  in subject  $j$  for cell type  $k$ , and  $s_j^k = \frac{\sum_{g \in G} \sum_{c \in C_j^k} \tilde{x}_{jgc}^k}{n_j^k}$  is the cell size defined as the average number of total mRNA molecules of cell type  $k$  in subject  $j$ . Let  $\mathbf{R}_j = [r_{jg}^k]_{G \times K}$  be the cell type specific reference matrix of subject  $j$ , then

$$r_{jg}^k = \frac{\sum_{c \in C_j^k} \tilde{x}_{jgc}^k}{n_j^k} = \theta_{jg}^k s_j^k.$$

To estimate the cell type proportions  $\mathbf{p}_j = (p_j^1, \dots, p_j^K)^T$  from (S.2) and address the heteroscedasticity (i.e., the unequal variance of  $\epsilon_{jg}$ ), we use the weighted least squares procedure. Namely,

$$\hat{\mathbf{p}}_j = \underset{\mathbf{p}_j \geq \mathbf{0}}{\operatorname{argmin}} \sum_{g \in G} w_{jg} \left( y_{jg} - n_j \sum_{k \in K} p_j^k r_{jg}^k \right)^2, \quad (\text{S.3})$$

where  $w_{jg}$  is a nonnegative weight. The nonnegative constraint on  $\mathbf{p}_j$  in (S.3) comes from the fact that cell type proportions are nonnegative. The reference  $r_{jg}^k$  in (S.3) is not directly observed and is typically replaced by an estimation  $\hat{r}_{jg}^k$  using the single-cell gene expression data, which we discuss in detail in the next subsection.

## 2.2 Technological discrepancies between sequencing platforms

Let  $x_{jgc}^k$  be the observed gene expression from the sc/snRNA-seq data, then we assume an additive model

$$x_{jgc}^k = \tilde{x}_{jgc}^k + \xi_{jgc}^k, \quad (\text{S.4})$$

where  $\xi_{jgc}^k$  is the measurement noise at the single-cell level, with mean zero and variance  $\delta_{gk}^2$ . The estimated cell type specific reference matrix  $\hat{\mathbf{R}}_j$  from the sc/snRNA-seq data is then expressed as

$$\hat{r}_{jg}^k = \frac{\sum_{c \in C_j^k} x_{jgc}^k}{n_j^k} = \theta_{jg}^k s_j^k + \frac{\sum_{c \in C_j^k} \xi_{jgc}^k}{n_j^k}. \quad (\text{S.5})$$

Then, by incorporating (S.5) into (S.2), we get

$$y_{jg} = n_j \sum_{k \in K} p_j^k (\hat{r}_{jg}^k - \frac{\sum_{c \in C_j^k} \xi_{jgc}^k}{n_j^k}) + \epsilon_{jg} = n_j \sum_{k \in K} p_j^k \hat{r}_{jg}^k - \gamma_{jg} + \epsilon_{jg},$$

where  $\gamma_{jg} = \sum_{k \in K} \sum_{c \in C_j^k} \xi_{jgc}^k$  is the accumulated measurement noise for gene  $g$  in subject  $j$  at the single-cell level. Then, the wNNLS framework for estimating the cell type proportions becomes

$$\begin{aligned} \hat{\mathbf{p}}_j &= \underset{\mathbf{p}_j \geq \mathbf{0}}{\operatorname{argmin}} \sum_{g \in G} w_{jg} \left( y_{jg} - n_j \sum_{k \in K} p_j^k \hat{r}_{jg}^k \right)^2 \\ &= \underset{\mathbf{p}_j \geq \mathbf{0}}{\operatorname{argmin}} \sum_{g \in G} w_{jg} \left[ \left( \tilde{y}_{jg} - n_j \sum_{k \in K} p_j^k \hat{r}_{jg}^k \right) + (\epsilon_{jg} - \gamma_{jg}) \right]^2. \end{aligned} \quad (\text{S.6})$$

Note that each residual  $\left(y_{jg} - n_j \sum_{k \in K} p_j^k \hat{r}_{jg}^k\right)$  in the wNNLS framework (S.6) is the summation of two terms, the true estimation error  $\left(\tilde{y}_{jg} - n_j \sum_{k \in K} p_j^k \hat{r}_{jg}^k\right)$  that we aim to minimize and the technological discrepancy  $(\epsilon_{jg} - \gamma_{jg})$  introduced by the two different sequencing platforms. Biologically, we expect both error terms  $\epsilon_{jg}$  and  $\gamma_{jg}$  are small so that the technological discrepancy is trivial and does not affect the estimation. In reality, however, either term could be large. Since both error terms are unlikely to be correlated, there is little chance of having two large terms that could be canceled out. In this scenario, the residual is dominated by the technological discrepancy rather than the true estimation error. Therefore, technological discrepancies can hinder the wNNLS framework from producing accurate estimates of cell type proportions.

### 2.3 A partitioned loss function in DeMixSC

The first innovation of DeMixSC is to employ a partitioned loss function to mitigate the negative effect of technological discrepancies. DeMixSC utilizes the differential expression (DE) analysis on matched bulk and pseudo-bulk data to categorize genes into two groups: technologically stable genes,  $G_1$ , and genes that are more susceptible to technological discrepancies,  $G_2$ . Genes from  $G_2$  hold larger values in technological discrepancies, which inadvertently steers the wNNLS framework to prioritize the minimization of the technological discrepancies over the true estimation error.

To address this issue, DeMixSC introduces a partitioned loss function, grounded in gene sets exhibiting distinct technological variances, and estimates the cell type proportions  $\mathbf{p}_j$  by

$$\hat{\mathbf{p}}_j = \underset{\mathbf{p}_j \geq \mathbf{0}}{\operatorname{argmin}} \left[ \sum_{g \in G_1} w_{jg} \left( y_{jg} - n_j \sum_{k \in K} p_j^k \hat{r}_{jg}^k \right)^2 + \sum_{g \in G_2} w_{jg} \left( \frac{y_{jg}}{a} - n_j \sum_{k \in K} p_j^k \frac{\hat{r}_{jg}^k}{a} \right)^2 \right], \quad (\text{S.7})$$

where  $a > 0$  serves as an adjustment for the expressions of genes within  $G_2$ , reducing the contribution of the squared residuals from  $G_2$  that are mainly dominated by technological discrepancies. As a result, we mitigate the adverse consequences stemming from the undesired technological discrepancies. The number of  $G_2$  genes and the adjustment  $a$  are user-defined tuning parameters for DeMixSC, and we set  $G_2$  to have 5000 genes and  $a$  to be 1000 by default after we examined the performance of DeMixSC with different parameter settings (Supplementary Fig. 16).

### 2.4 A new weight function in DeMixSC

In the wNNLS framework (S.6), a typical choice of the weight function is the inverse of the variance of the measurement noises, which is unknown in our deconvolution problem. A common solution to resolve this issue is to estimate the weights iteratively, see [1, 2, 3, 4].

The second enhancement introduced in DeMixSC is a new weight function for the updated wNNLS model (S.7), which is defined as follows:

$$\begin{aligned} w_{jg}^{*-1} &= \hat{y}_{jg}^2 + (y_{jg} - \hat{y}_{jg})^2 + c \\ &= \left( n_j \sum_{k \in K} \hat{p}_j^k \hat{r}_{jg}^k \right)^2 + \left( y_{jg} - n_j \sum_{k \in K} \hat{p}_j^k \hat{r}_{jg}^k \right)^2 + c. \end{aligned} \quad (\text{S.9})$$

The weight function contains three terms: a squared fitting term within the bulk RNA-seq data, a squared residual, and a baseline constant  $c$ . The squared fitting in the first term in equation (S.9) accounts for genes with high expression levels in the sequencing data, to downweight these genes' impact on the loss function [3]. The squared residual in the second term in (S.9) is an estimate of the variance component, which accounts for the remaining variance in bulk RNA-seq data after fitting. It reduces the contribution of genes with large residuals and improves the fit to the data that is already well-estimated. The baseline constant  $c$  in (S.9) is for constraining the range of weights. It keeps the weights of all genes within a reasonable range to prevent extreme values in weights resulting in potential numerical issues. We set the constant  $c$  to 2 by default based on our empirical test (see Supplementary Fig 17).

### 3 Computation of the DeMixSC estimate

Following the existing wNNLS-based deconvolution method, DeMixSC computes its estimate through an iterative algorithm. We initiate the DeMixSC computational framework with a nonnegative ordinary least squares estimate:

$$\hat{p}_j^{(0)} = \underset{\mathbf{p}_j \geq \mathbf{0}}{\operatorname{argmin}} \left[ \sum_{g \in G_1} \left( y_{jg} - n_j \sum_{k \in K} p_j^k \hat{r}_{jg}^k \right)^2 + \sum_{g \in G_2} \left( \frac{y_{jg}}{a_g} - n_j \sum_{k \in K} p_j^k \frac{\hat{r}_{jg}^k}{a_g} \right)^2 \right]. \quad (\text{S.10})$$

Then, with the cell type proportion estimates from (S.10), we update the gene weights  $w_{jg}$  and re-estimate the cell type proportion by iteratively solving:

$$\hat{p}_j^{(l)} = \underset{\mathbf{p}_j \geq \mathbf{0}}{\operatorname{argmin}} \left[ \sum_{g \in G_1} w_{jg}^{*(l)} \left( y_{jg} - n_j \sum_{k \in K} p_j^k \hat{r}_{jg}^k \right)^2 + \sum_{g \in G_2} w_{jg}^{*(l)} \left( \frac{y_{jg}}{a_g} - n_j \sum_{k \in K} p_j^k \frac{\hat{r}_{jg}^k}{a_g} \right)^2 \right], \text{ for } l = 1, 2, \dots,$$

where  $w_{jg}^{*(l)}$  is updated using the previous estimate  $\hat{p}_j^{(l-1)}$ , i.e.,

$$\left[ w_{jg}^{*(l)} \right]^{-1} = \left( n_j \sum_{k \in K} \hat{p}_j^{k(l-1)} \hat{r}_{jg}^k \right)^2 + \left( y_{jg} - n_j \sum_{k \in K} \hat{p}_j^{k(l-1)} \hat{r}_{jg}^k \right)^2 + 2.$$

The final DeMixSC estimate is obtained when either the algorithm converges or reaches the prespecified maximum number of iterations.

## References

- [1] M. Dong, A. Thennavan, E. Urrutia, Y. Li, C. M. Perou, F. Zou, and Y. Jiang. SCDC: bulk gene expression deconvolution by multiple single-cell RNA sequencing references. *Briefings in bioinformatics*, 22(1):416–427, 2021.
- [2] J. Fan, Y. Lyu, Q. Zhang, X. Wang, M. Li, and R. Xiao. Music2: cell-type deconvolution for multi-condition bulk rna-seq data. *Briefings in Bioinformatics*, 23(6):bbac430, 2022.
- [3] D. Tsoucas, R. Dong, H. Chen, Q. Zhu, G. Guo, and G.-C. Yuan. Accurate estimation of cell-type composition from gene expression data. *Nature communications*, 10(1):2975, 2019.
- [4] X. Wang, J. Park, K. Susztak, N. R. Zhang, and M. Li. Bulk tissue cell type deconvolution with multi-subject single-cell expression reference. *Nature communications*, 10(1):380, 2019.