

EMT Multi-Omics

Deconvolution of proteomics data using scRNA-Seq

Ahmed Youssef

September 24, 2021

Contents

Introduction	1
Experiment summary	2
Approach	2
Data summary - Proteomics	2
Data summary - scRNA-Seq	3
Data summary - Bulk mRNA	3
Protein overlap	4
Pseudo-bulk RNA vs bulk RNA	4
Relationship between scRNA and bulk mRNA	4
Identify cell clusters	5
Changes in cell state composition over time	6
Re-construct pseudo-bulk data from single-cell data	6
Re-construct single-cell data from pseudo-bulk data	7
Variation of gene expression within cell type across timepoints	8
Constraining the solution to be positive	8
Predicting single-cell profiles using ridge regression	10

Introduction

The fundamental unit of all living organisms is the cell, and recent technological advances have granted us unprecedented opportunities to study life at this principal level. Proteins, through their networks of interactions, carry out most of the vital biological processes governing cellular functions, yet remain largely

unexplored in the single-cell space, representing crucial gaps in our knowledge of cell biology. While single-cell proteomics methods are still in their infancy, single-cell RNA sequencing (scRNA-Seq) has emerged in recent years as a powerful technology for defining cell states on a large scale, enabling breakthroughs in many areas of cell biology research, and begging the question of whether it can be used for making inferences at the protein level. **In this report, I explore the deconvolution of bulk proteomics data to the single-cell level using scRNA-Seq data.**

Experiment summary

Epithelial-to-mesenchymal transition (EMT) is a biological process in which epithelial cells gradually lose their adhesion and transition into mesenchymal cells. As one of the hallmarks of cancer progression, it is one of the long-standing interests of the biomedical research community. Towards profiling this process, protein and RNA samples were extracted from cells at 8 different timepoints during EMT and multiple layers of omics data were generated. These omics layers include proteomics, transcriptomics, phosphoproteomics, secretome, exosome among others. A pre-print with more details on the experiment and generated data can be found on bioRxiv [here](#) (Paul et al, 2021). This report is interested in the scRNA-Seq, microarray, and proteomics datasets generated in this study.

Approach

Bulk proteomics data gives a view of the aggregated protein abundance from all cell types within a sequenced sample. Using single-cell data, derived from the same samples, we can investigate the sample heterogeneity by estimating proportions of cell types within a bulk sample. We cannot reliably use these proportions to directly estimate the contribution of each population to each gene/protein's expression at the bulk-level however, since there is low correlation between RNA and protein levels of the same genes due to multiple biological factors such as alternative splicing and post-translational modifications. Leveraging the timepoints present in this dataset, which conveniently show shifts in cell type abundances across time, we can instead look for changes in cell-type proportions and corresponding changes in bulk-level protein abundance as suggestive of relationships between specific cell types and specific proteins. This information can then potentially be used to estimate the contribution of individual cell types to the bulk proteomics measurements.

The step-by-step process is outlined below:

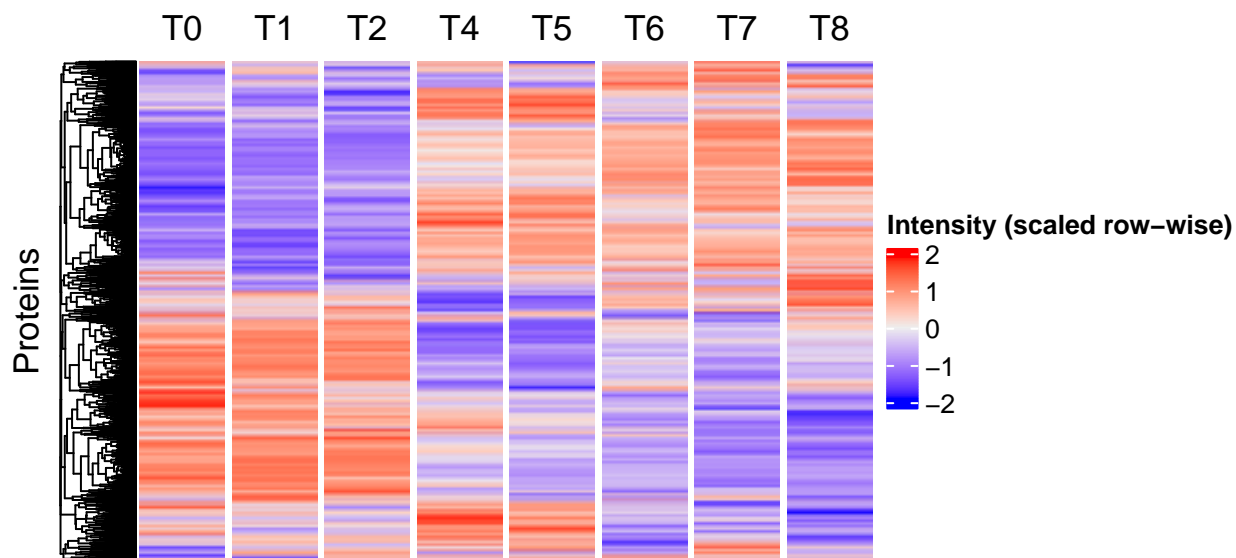
- 1) Pool together all single cells across timepoints into one scRNA-Seq dataset
- 2) Cluster the pooled scRNA data to identify cell states
- 3) Compute proportion of each cell state in each timepoint using the cluster labels derived from the pooled data
- 4) For each gene-cell state pairing, compute the Pearson correlation coefficient between the 8 timepoint proteomics measurements (gene) and the cell cluster proportions (scRNA)

Data summary - Proteomics

The bulk proteomics data was generated in the Emili Lab using standard mass-spectrometry. Summary of the dataset follows:

- 6,426 proteins
- 10 different timepoints
- Three replicates

The average intensity across replicates was computed for each protein in each timepoint. Timepoints 3 and 9 were removed since they are not present in the scRNA data.



Data summary - scRNA-Seq

The bulk proteomics data was generated in the Emili Lab using standard mass-spectrometry. Summary of the dataset follows:

- 9,785 genes
- 1,913 cells (~200 cells per timepoint)
- 8 different timepoints

Prior to this summation, genes with zero variance as well as those with non-zero counts in less than 5% of all cells were removed. This removed 17 genes (0.2% of all genes).

Table 1: Number of cells in each timepoint

T0	T1	T2	T4	T5	T6	T7	T8
367	208	246	190	247	217	236	202

Data summary - Bulk mRNA

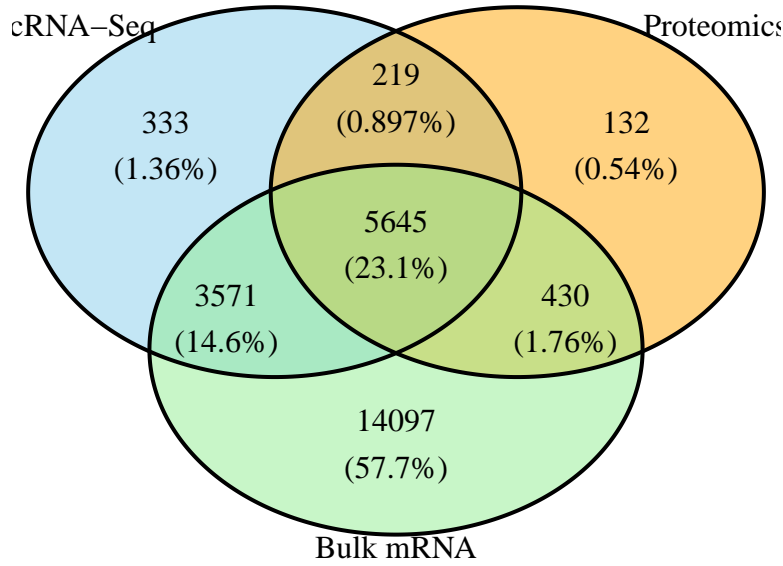
The bulk mRNA data comes from a microarray experiment. Summary of the dataset follows:

- 23,743 genes
- 10 different timepoints
- Three replicates

The average intensity across replicates was computed for each protein in each timepoint. Timepoints 3 and 9 were removed since they are not present in the scRNA data.

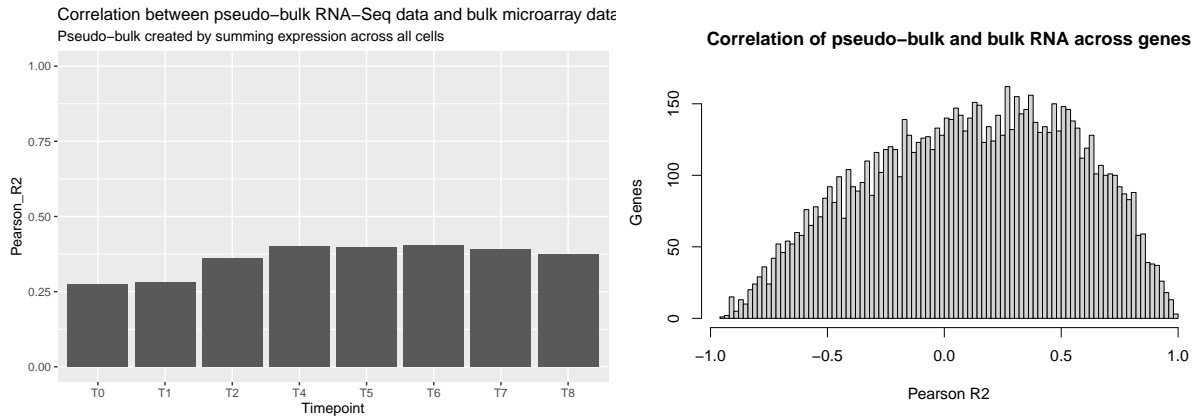
Protein overlap

The venn diagram below shows the overlap of the identified proteins in the datasets.



Pseudo-bulk RNA vs bulk RNA

A pseudo-bulk RNA dataset is created for each timepoint by summing the gene counts of all cells within the timepoint. This pseudo-bulk data is then compared to the actual microarray bulk mRNA data present for each timepoint. The cross-timepoint measurements for each genes were also correlated. The distributions of these correlations are showed below.



Relationship between scRNA and bulk mRNA

To investigate the ability to recover the scRNA data from the pseudo-bulk data, we start by clustering the scRNA data to identify the cell clusters followed by creating the following matrices:

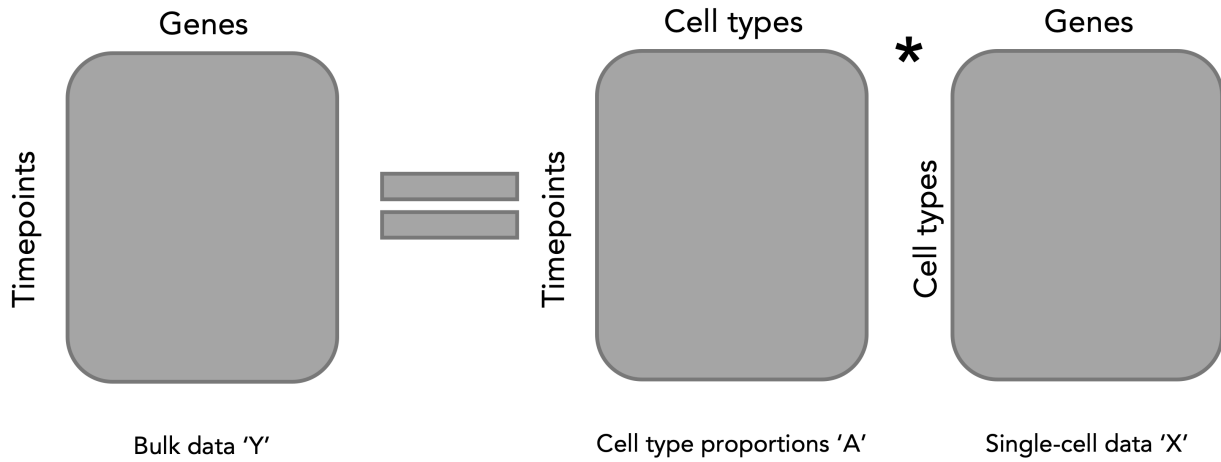
- Matrix A of dimensions $timepoints \times clusters$. (cell type proportions)

- Matrix X of dimensions $clusters \times genes$. (cluster-averaged single-cell RNA data)
- Matrix Y of dimensions $timepoints \times genes$. (pseudo-bulk RNA data)

We then attempt to re-create the single-cell matrix X data by computing $Y = AX'$.

Bulk deconvolution

$$Y = AX$$

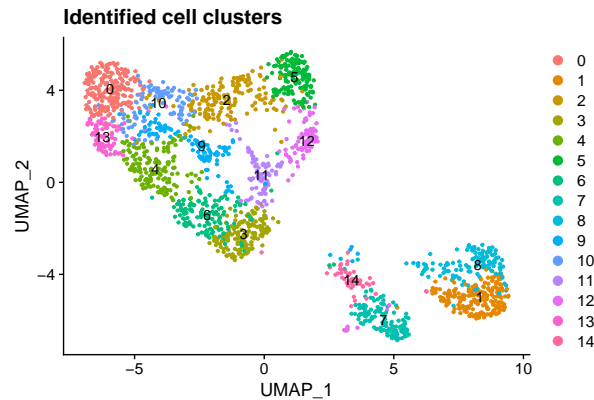


Identify cell clusters

The cell states in our dataset are identified in an unsupervised manner based on similarity of gene expression profiles. All cells from all timepoints are pooled together for this analysis.

For data pre-processing, we remove the genes with low expression counts, retaining genes with a minimum of 3 counts in at least 3 cells. This removed 1,240 genes (13% of all genes). On average, each cell expressed ~3,600 genes after processing. [Seurat](#) is then used to cluster the cells based on the 2,000 most variable genes.

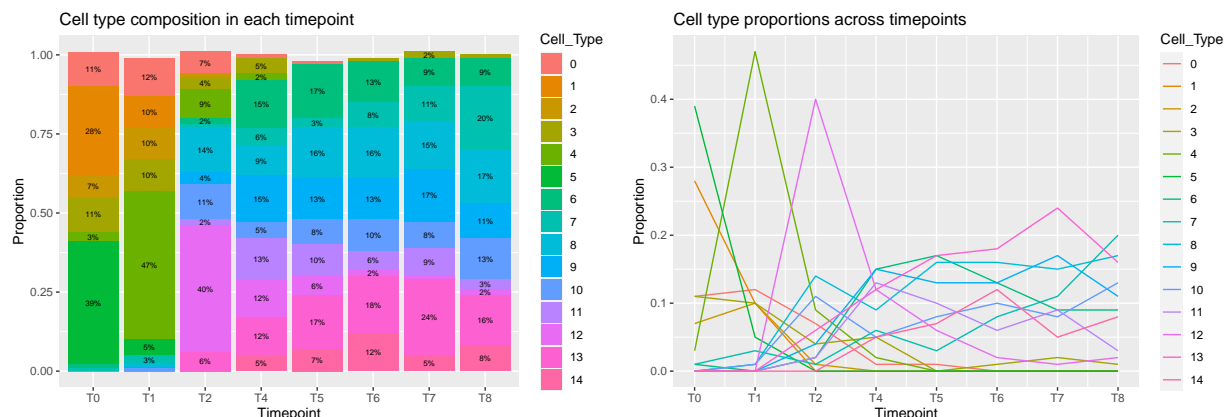
The initial model defines 7 cell clusters at a Seurat resolution of 0.4. We intentionally begin with a number of clusters less than the number of timepoints to avoid creating an undetermined problem. The below UMAP plot visualizes the identified cell clusters.



Changes in cell state composition over time

Since this dataset is investigating cells undergoing EMT, the cell population abundances are changing over time. The below figures visualize these proportion changes.

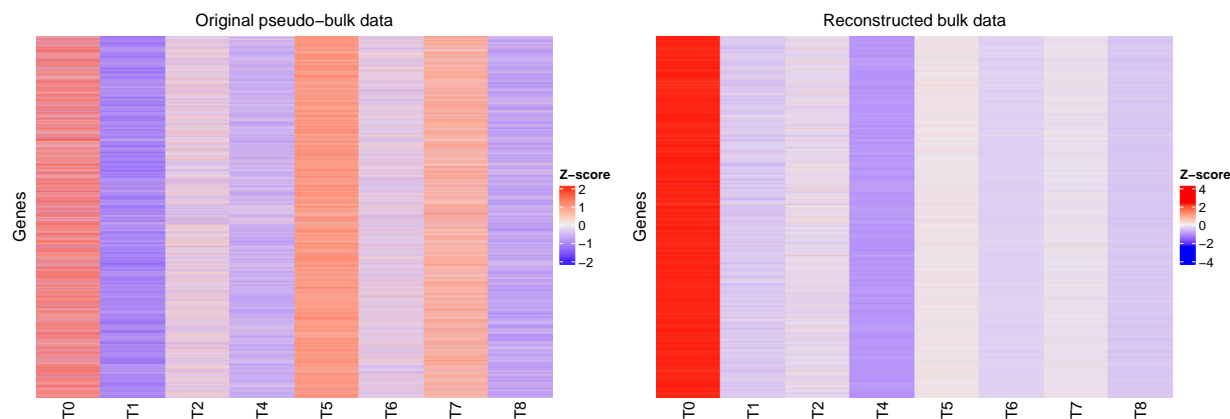
Note: this approach assumes the single-cell data accurately captures the sample heterogeneity. In practice, biased cell sampling upstream could lead to an inaccurate view of sample heterogeneity here.

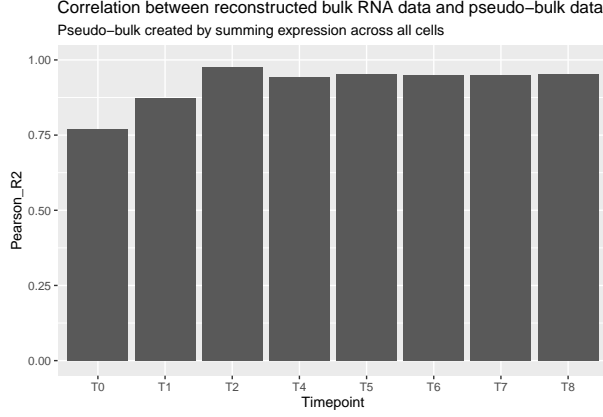


Re-construct pseudo-bulk data from single-cell data

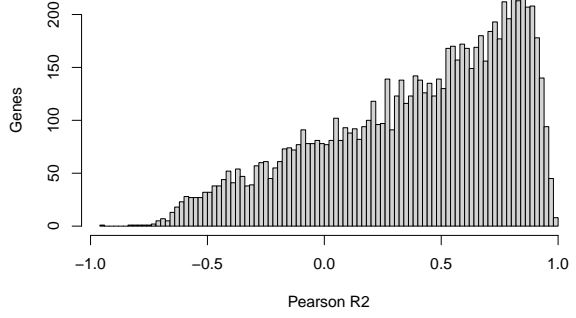
We now attempt to re-create the pseudo-bulk data (Y) from the single-cell data using the simple formula $Y' = AX$. The matrix X was constructed by averaging the gene expression of each cluster. The matrix A contains the number of cells from each cell cluster in each timepoint.

The below heatmaps show the pseudo-bulk data and the re-constructed data side-by-side, with the 8,528 genes in the same order. The heatmaps are scaled by row for visualization purposes. This reconstructed bulk data is then compared to the pseudo-bulk data in each timepoint. Correlations were also computed across the cross-timepoint measurements for each gene. The distributions of these correlations are shown below.

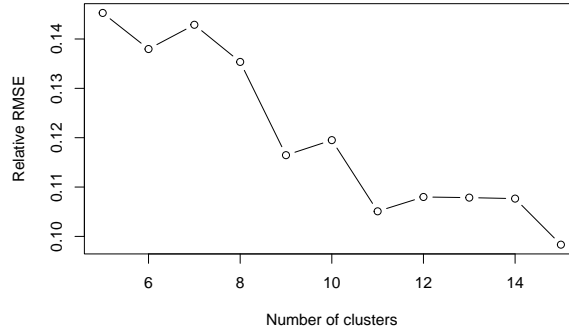




Correlation of reconstructed bulk RNA and pseudo-bulk RNA across g



Relative RMSE between Y and Y' for different no. of clusters

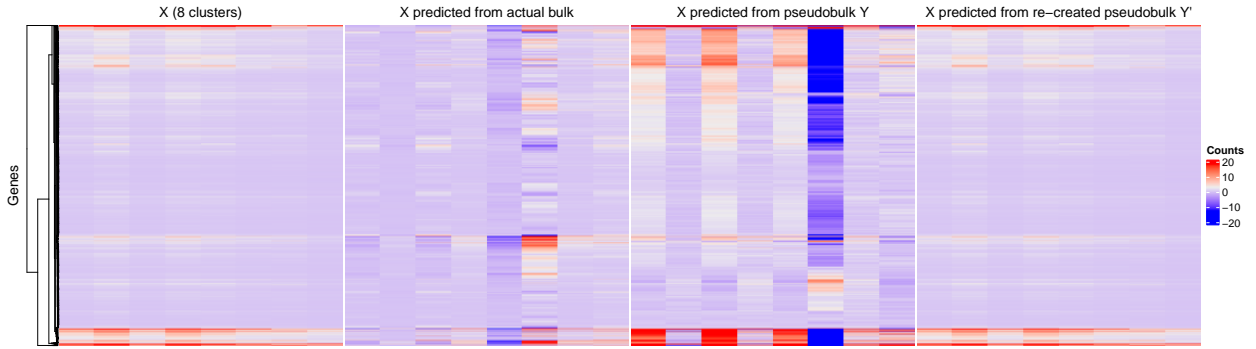
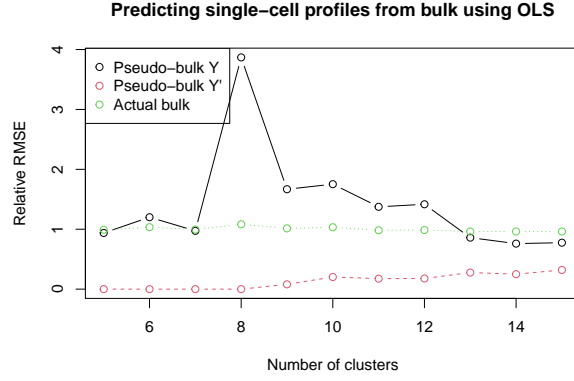


Re-construct single-cell data from pseudo-bulk data

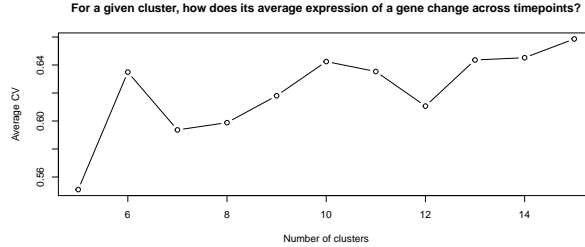
In this section, we attempt to re-create the single-cell data from the pseudo-bulk data and the timepoint-specific cell cluster counts. Based on the formula $Y = AX$ outlined in previous sections, given the pseudobulk matrix Y and the timepoint-specific cell counts ‘mixing’ matrix A , we aim to compute X using the formula $X' = (A^T A)^{-1} (A^T Y)$, which is essentially the pseudo-inverse of A multiplied by the pseudo-bulk Y . Recall that the pseudo-bulk is computed by summing up the counts of individual cells in each timepoint.

We vary the number of cell clusters by varying the *resolution* parameter in Seurat’s clustering algorithm. We try this method for a number of clusters ranging between 5-15. For each given number of clusters, we solve the formula above to predict X' . We then compare this to the actual observed X from the scRNA-Seq data, which is the per-cluster average of gene counts. The errors are reported as relative RMSE using the formula $\|X - X'\| / \|X\|$, where $\|X\|$ is the Frobenius norm of the matrix (the square root of the sum of squares of the matrix elements).

This process above is repeated once using Y' as the bulk data, where $Y' = AX$ is computed first prior to computing X' as above, and once using the actual bulk data in this dataset. The resultant cluster-RMSE relationships are shown in the below plot.



Variation of gene expression within cell type across timepoints



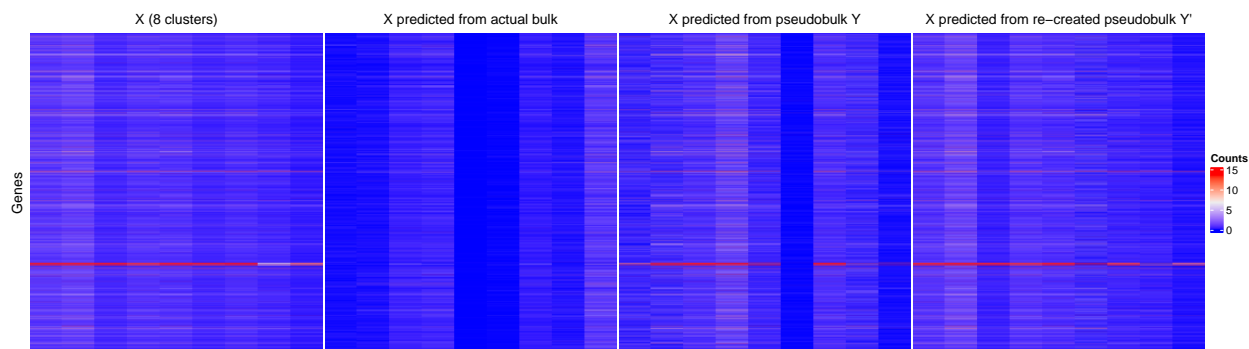
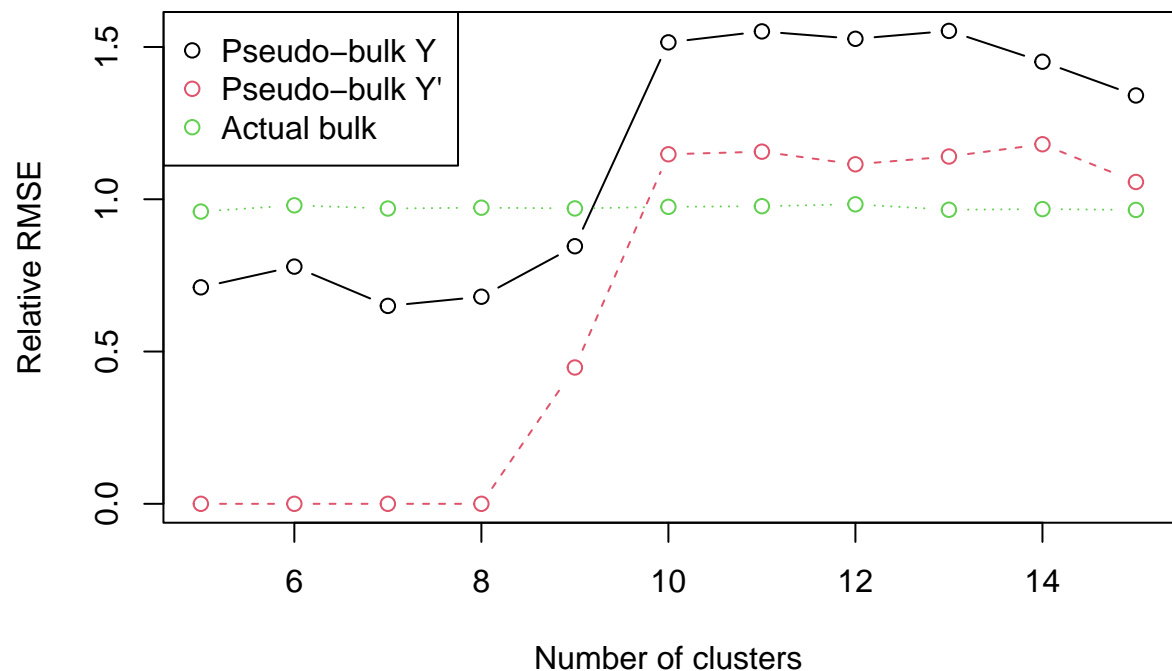
Constraining the solution to be positive

The (Tsoucas et al)[<https://www.nature.com/articles/s41467-019-10802-z>] publication proposed a strategy for RNA deconvolution based on the OLS formulation we used in this report, with the addition of constraining the solution to be positive since there are no possible negative gene counts, and increasing the weight of cell types with low average expression levels. Since we typically have more cell types than time points, our problem is underdetermined and as such there are multiple possible solutions. In this section, we follow their strategy for constraining the solution to be non-negative.

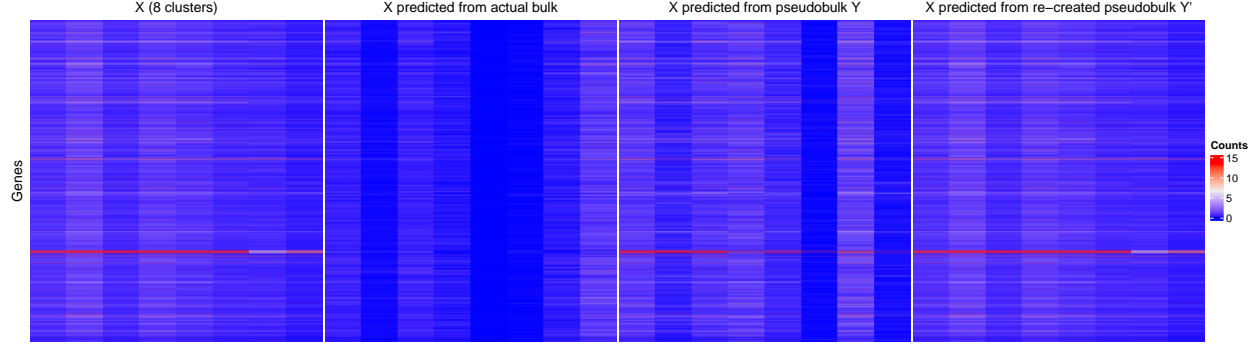
We solve $\hat{X} = \min_{x \geq 0} \|AX - Y\|^2$ which algebraically is the same as solving $\hat{X} = \min_{x \geq 0} (-2Y^T A X + X^T A^T A X)$. The R function (solve.QP)[<https://www.rdocumentation.org/packages/quadprog/versions/1.5-8/topics/solve.QP>] from the *quadprog* package is used to solve this equation one gene at a time to estimate the expression profile at the cell-type level. This algorithm, however, expects a positive definite matrix. The matrix $A^T A$ is positive definite when A has more rows than columns – i.e., when the problem is a

least-squares problem. When we have more clusters than timepoints then A has more columns than rows, and so the matrix is not invertible, and so it is not positive definite. As a workaround, we use the function *nearPD()* to find the nearest positive matrix instead.

Solving constrained OLS using nearest positive-definite matrix



As an alternative to the *nearPD* function, we add a small scaling of the identity matrix (10^{-9}) to $A^T A$, which will make it positive definite and allow *solve.QP* to work. The result is shown below.



Predicting single-cell profiles using ridge regression

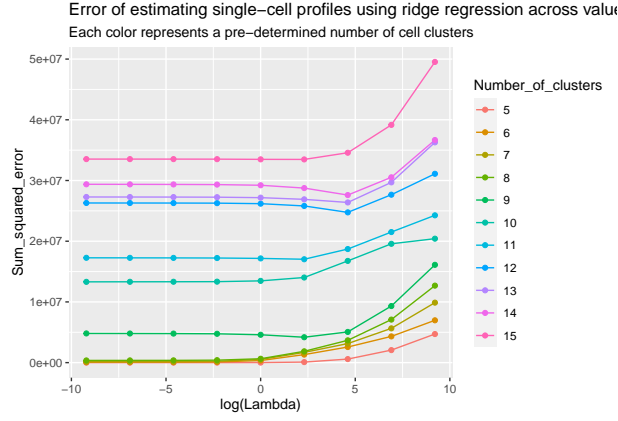
The condition number of $A^T A$ for the case where we have 8 single-cell clusters, i.e. as many clusters as there are timepoints, is noticeably high (15,828) compared to a smaller number of clusters which suggests that it is a near-singular matrix. This would explain the relatively high error in estimating the single-cell profiles. In an attempt to address this, we use ridge regression to estimate the single-cell profiles instead.

Ridge regression is similar to linear regression but with the addition of a regularization term λ . A range of values for λ are tested to decide on the optimal one. The R package *glmnet* is used for this task. We solve the problem by fitting a ridge regression model to each gene. The full procedure is outlined below:

Input: Number of clusters N , Single-cell expression matrix SC

Procedure:

1. Cluster single-cell matrix SC into N clusters using *Seurat*
2. Construct cluster * gene matrix \mathbf{X} by computing the average expression of each gene in each cluster
3. Construct timepoint * cluster mixing matrix \mathbf{A} by counting the number of cells in each cluster in each timepoint
4. Construct timepoint * gene pseudobulk matrix \mathbf{Y} using the formula $Y = AX$
5. Define $\text{LambdaSet} = \{10^0, 10^1, 10^2, 10^3, 10^4, 10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}\}$
6. For each column of Y , denoted as y :
 7. Fit ridge regression model $Beta$ using the function `Beta = glmnet(x = A, y = y, lambda = LambdaSet, alpha = 0)`
8. For each column in X , denoted as x :
 - For each lambda L :
 1. Compute $x' = \text{column } L \text{ of } Beta$
 2. Compute error in predicting x as $\|x - x'\|^2$
 3. Accumulate error for each lambda as the sum of errors across genes
 4. Select lambda with minimum accumulated error



For each number of clusters, we select the lambda which led to the minimal sum squared error when comparing the predicted single-cell profiles to the actual ones. Those selected lambdas are shown below.

Table 2: Ridge regression - Lambda values that minimize error

Number_of_clusters	best_lambda
5	1e-04
6	1e-04
7	1e-04
8	1e-04
9	1e+01
10	1e-04
11	1e+01
12	1e+02
13	1e+02
14	1e+02
15	1e+01

Comparing ridge regression to OLS

