# UROP Notes

MICHAEL DIAO

## 1 RCTD

Our ultimate goal is to determine the fractional contributions of each cell type to a particular sample. We do this by maximum likelihood, using the following hierarchical model:

$$Y_{i,j} \mid \lambda_{i,j} \sim \text{Poisson}(N_i \lambda_{i,j})$$
$$\log \lambda_{i,j} = \log(\vec{\beta}_i \cdot \vec{\mu}_j) + \alpha_i + \gamma_j + \varepsilon_{i,j},$$

where

- $Y_{i,j}$ is the random variable corresponding to the observed expression of gene $j$ at pixel $i$,

- $N_i$ is the number of transcripts for pixel $i$,

- $\vec{\beta}_i$ is the $K$-dimensional row vector of contributions from each cell type (where $K$ is the number of cell types in question) at pixel $i$,

- $\vec{\mu}_j$ is the $K$-dimensional column vector of mean expressions of gene $j$ for each cell type,

- $\alpha_i$ is a fixed pixel-specific effect.

- $\gamma_j$ and $\varepsilon_{i,j}$ are random effects that introduce noise. $\gamma_j$ in particular is intended to account for platform effects that may over- or underrepresent certain genes. We let these be normally distributed with mean 0 and variance $\sigma_\gamma$, $\sigma_\varepsilon$ respectively.

Therefore, determining the fractional contributions of each cell type reduces to finding the maximum likelihood parameter $\vec{\beta}_i$ for each $i$.

> **Question 1.1.** Now we have a ton of parameters, potentially thousands. $\beta$ alone introduces $K \times J$ of them. How do we do any useful estimation here?

We proceed in the following steps:

1. **Supervised estimation of cell type profiles**.

   Using a reference dataset, we estimate the parameters $\mu_j$ of expression levels for gene $j$, giving $\hat{\vec{\mu}}_j$ which will be used in the next steps.

   We can do this by obtaining a (e.g. scRNA-seq) reference annotated with cell types, after which $\vec{\mu}_j$ can be estimated as the empirical average normalized expression of gene $j$ within each cell type.

2. **Gene filtering**.

   Using the estimated expression profiles $\hat{\vec{\mu}}_j$, we filter out genes that are not highly variable across cell types.

   We can do this by taking the expression profiles $\hat{\vec{\mu}}_j$ and selecting genes with a minimum average expression and sufficiently high variance.

   > does this discard genes that are, say, only expressed in one cell type? (average might be low, but gene could be good marker)

3. **Platform Effect Normalization**.

   With an estimate for $\vec{\mu}_j$, it turns out we now have a way to estimate the platform effects $\gamma_j$ as well. The idea is that we can consider the average observed expression across pixels

   $$M_j = \frac{1}{I} \sum_{i=1}^{I} Y_{i,j},$$

   whence

   $$\log \mathbb{E}_{M_j|\vec{\lambda}_j}\left[M_j \mid \lambda_{1,j}, \ldots, \lambda_{I,j}\right] = \log\left(\frac{1}{I} \sum_{i=1}^{I} N_i \lambda_{i,j}\right)$$

   $$= \log\left(\frac{1}{I} \sum_{i=1}^{I} N_i \exp\left(\log(\vec{\beta}_i \cdot \vec{\mu}_j) + \alpha_i + \gamma_j + \varepsilon_{i,j}\right)\right)$$

   $$= \gamma_j + \log\left(\frac{1}{I} \sum_{i=1}^{I} N_i (\vec{\beta}_i \cdot \vec{\mu}_j) \exp\left(\alpha_i + \varepsilon_{i,j}\right)\right)$$

   $$= \gamma_j + \log\left(\frac{1}{I} \sum_{i=1}^{I} \left(\sum_{k=1}^{K} \beta_{i,k} \cdot \mu_{k,j}\right) N_i \exp\left(\alpha_i + \varepsilon_{i,j}\right)\right)$$

   ¡++¿

¡++¿