

Molecular Classification 2

BIOF2014

Problem

We want to classify tumours based on their molecular characteristics, as determined by RNA expression profiles. Each expression profile consists of the expression levels of selected genes, measured as counts of molecules.

We have previously mathematically derived posterior predictive distribution $p(\tilde{y} | \tilde{\mathbf{x}}, \mathbf{X}, \mathbf{y})$. for a new class label \tilde{y} given new data $\tilde{\mathbf{x}}$.

If we are willing to tolerate computational inefficiency, we can also use an Markov chain Monte Carlo sampler such as [Stan](#) to estimate the posterior predictive distribution (or any quantities of interest).

Model

Given training data with expression profiles $\mathbf{X} = [\mathbf{x}_i^\top]$ and class labels $\mathbf{y} = [y_i]$ for $i \in \{1 \dots N\}$, we want to predict the unknown label \tilde{y} of a new sample with expression profile $\tilde{\mathbf{x}}$.

Each expression profile consists the detected counts of transcript molecules of J genes, so $\mathbf{x}_i, \tilde{\mathbf{x}} \in \mathcal{N}_0^J$. Suppose there are K classes, so $y_i, \tilde{y} \in \{1 \dots K\}$.

Define our model as follows:

$$\begin{aligned}\mathbf{X}_i = \mathbf{x} \mid Y_i = y, \boldsymbol{\eta}_y &\sim \text{Multinomial}(m_i \mid \boldsymbol{\eta}_y) \\ \boldsymbol{\eta}_k &\sim \text{Dirichlet}(\mathbf{c}_k) \\ Y_i &\sim \boldsymbol{\theta} \\ \boldsymbol{\theta} &\sim \text{Dirichlet}(\mathbf{d}),\end{aligned}$$

where $\mathbf{c}_k \in \mathcal{R}_{\geq 0}^J$ and $\mathbf{d} \in \mathcal{R}_{\geq 0}^K$ are hyperparameters.

Tasks

1. Implement the classification model in Stan and R and sample from $p(\tilde{y} | \tilde{\mathbf{x}}, \mathbf{X}, \mathbf{y})$.
2. Implement the previously derived $p(\tilde{y} | \tilde{\mathbf{x}}, \mathbf{X}, \mathbf{y})$.
3. Compare the empirical pmf of from MCMC sampling vs. the mathematically-derived posterior predictive distribution $p(\tilde{y} | \tilde{\mathbf{x}}, \mathbf{X}, \mathbf{y})$.