

Multiple sample comparison in flow cytometry data with flowMap

Chiaowen Joyce Hsiao

Center for Bioinformatics and Computational Biology
University of Maryland, College Park

`chsiao@umiacs.umd.edu`

Modified: September 24th, 2013. Compiled: July 13, 2014

Contents

1	Introduction	2
2	Overview of the algorithm	2
3	Data preparation	3

1 Introduction

Flow cytometry (FCM) is a powerful single-cell technology that provides high-throughput data on cellular features, such as size, complexity, and antibodies. The data is processed one sample at a time. Homogeneous groups of cells, also known as cell populations, are then identified via automatic or manual gating methods. These cell populations may vary in proportion, shape or levels of a particular antibody markers between samples. A major step in the downstream analysis of flow cytometry data is to analyze cell population differences for phenotype comparisons. `flowMap` implements a nonparametric variate test to compare cell populations between samples. Our method can accommodate the high-dimensional features of FCM data and also the sample variation in distributions. Details of the method is described in [1].

2 Overview of the algorithm

`flowMap` implements the nonparametric Friedman-Rafsky (FR) multivariate run test to compare and match cell populations between samples. p-values of the FR statistic are calculated for each population comparison. Populations are considered a match if the p-values are above a Bonferroni-corrected cutoff. Moreover, we employed two approaches to calculate the p-values: 1) finding the percentile of the observed FR statistics in a standard normal distribution, and 2) finding the percentile of the statistics in a empirical null distribution of the FR statistics. The former assumes that the FR statistics follows a standard normal distribution, while the latter takes no distributional assumptions. Hereafter, the former p-value is referred to as the *theoretical p-value*, and the latter is referred to as the *empirical p-value*. Both p-values are provided in the output.

The goal of our algorithm is to identify matched versus mismatched cell populations by comparing each FCM sample to a selected reference sample. Thus, two cell populations matched to the same reference are considered to be similar to each other as well. This method reduces the computational complexity. In addition, the algorithm allows the user to choose a reference sample for mapping or to construct a reference sample from the FCM test samples. The general flow of the diagram is as follows:

Denote S_o as the reference sample with m populations, and S_i as the test sample with n_i populations where $i = 1, \dots, n$. Then,

1. Compare every S_i with S_o ,

Step 1: compute FR statistics of the $n_i \times m$ population pairs,

Step 2: compute p-values for the FR statistics,

Step 3: identify a population pair as matched if the p-value is less than $0.01/(m \times n_i)$,

2. Reassign S_i population labels to the matched S_o population label or a new unique population label if there is no match in S_o .
3. Make a metaset of cell population labels by combining the matched and mismatched populations across all samples S_1, \dots, S_n .

3 Data preparation

Here we assume that the data have been normalized and transformed according to appropriate flow cytometry data procedure. The input data can be in txt format or as data.frames, where the rows are the event (cell) data. The columns are consisted of the features and also the cell population identifying number as the last column of the data. Below is an example data `Sample1`. There are 25,809 events in total with 5 feature markers (CD20,CD24,CD27,CD38,IgD). The last column of the data indexes the cell population labels.

```
list()
```