

Introduction to pre-processing – lab

Martin Morgan

28 January 2010

This lab introduces pre-processing and quality assessment, focusing on Affymetrix single channel data. It is based on Chapter 3 of *Bioconductor Case Studies*, by Hahne, Huber, Gentleman, and Falcon.

1. Briefly read through the chapter 3 handout.
2. Ensure that you have the CLL package installed. To do this, start an R session and type the command

```
> library(CLL)
```

If R responds with `Error in library(CLL) : there is no package called 'CLL'` then use the ethernet connection to install the package:

```
> source("http://wilson2.fhcrc.org/installScripts/biocLite.R")
> biocLite(CLL)
```

(use `source("http://bioconductor.org/biocLite.R")` for access after the course is complete.)

3. (Optional) Copy the CEL files from the thumb drive to your disk.

Exercise 1

Attach the *CLL* package, and use the `data` function to load *CLLbatch* into your work space. The *CLLbatch* object contains raw probe values; they have not been summarized or normalized in any way. What is the class of *CLLbatch*? Can you find the help page describing this class? What about a help page describing the *CLLbatch* itself? From looking at what is printed when you type *CLLbatch* at the command line, what basic information (e.g., number of genes, number of samples, the chip on which the expression data were collected) can you obtain about *CLLbatch*?

Exercise 2

As an exercise, we'll add some phenotypic data *CLLbatch*. Do this by using `data(disease)`. Take a minute to explore this object. What class is it, how many rows does it have, what are its column names?

Update the row names of the phenotype `data.frame` to reflect the `SampleID`, and remove the `.CEL` file extension from the `sampleNames` of the *AffyBatch*.

Create a vector `mt` that tells us how the row names of the phenotype `data.frame` align with the column names of the `AffyBatch`.

Create the ‘metadata’ describing the phenotype data, then an `Annotated-DataFrame` combining the phenotype data and metadata, and finally add this to the `CLLbatch AffyBatch`. Use `mt` to make sure the order of the phenotype and expression data is consistent.

One array is missing disease status. Drop this from further analysis.

Exercise 3

There is an array quality metrics report for `CLLbatch`; we’ll review it in class. The report suggests that array `CLL1` is suspect, so remove it from further analysis.

Exercise 4

Pre-process the `AffyBatch` using `rma`. What is the class of the resulting object? How many probesets are there in this object? Can you extract the summarized expression values (using the `exprs` function)? What’s happened to the phenotypic data that we added to `CLLbatch`?

Often `just.rma` is used to go directly from `CEL` file to `ExpressionSet` in a fast and memory efficient way.

Exercise 5

Functions can often help to reduce the need to repeatedly perform routine tasks, and in the process make work flow less error prone. The `NUSE` boxplot in the chapter represents two steps: processing an `AffyBatch` object using `fitPLM`, and then plotting the results. Capture these steps in a function, and add it to your package.

Some steps might include:

1. Decide on what your function will do: for an `AffyBatch` object, fit a probe-wise linear model. Then create a `NUSE` plot.
2. Decide what the inputs and outputs of your function must be. The input will be an `AffyBatch` object. The ‘output’ is really a side effect, the creation of a plot.
3. Create a simple function taking the appropriate inputs and outputs, but doing nothing else. Choose a function name and arguments with an eye toward making their purpose clear. Create the function in a file in the `R` directory of your simple package.
4. Fill in the function with the details.
5. Test the function.