

Normalization: Bioconductor's marray package

Yee Hwa Yang¹ and Sandrine Dudoit²

April 25, 2007

1. Department of Medicine, University of California, San Francisco, jean@biostat.berkeley.edu
2. Division of Biostatistics, University of California, Berkeley,
<http://www.stat.berkeley.edu/~sandrine>

Contents

1	Overview	1
2	Getting started	2
3	Normalization using robust local regression	2
4	Normalization functions	3
4.1	Simple normalization function <code>maNorm</code>	3
4.2	Simple scale normalization function <code>maNormScale</code>	4
4.3	General normalization function <code>maNormMain</code>	4
5	Normalization of Swirl zebrafish microarray data	5
5.1	Using simple function <code>maNorm</code>	5
5.2	Using simple function <code>maNormScale</code>	6
5.3	Using main function <code>maNormMain</code>	7
5.4	Plots	7

1 Overview

This document provides a tutorial for the normalization component of the `marray` package. Greater details on the packages are given in [Dudoit and Yang(2002)]. This package implements robust adaptive location and scale normalization procedures, which correct for different types of dye biases (e.g. intensity, spatial, plate biases) and allow the use of control sequences spotted onto the array and possibly spiked into the mRNA samples. Normalization is needed to ensure that observed differences in intensities are indeed due to differential expression and not experimental artifacts; fluorescence intensities should therefore be normalized before any analysis which involves comparisons among genes within or between arrays.

2 Getting started

To load the `marray` package in your R session, type `library(marray)`. We demonstrate the functionality of this R packages using gene expression data from the Swirl zebrafish experiment. These data are included as part of the package, hence you will also need to install this package. To load the swirl dataset, use `data(swirl)`, and to view a description of the experiments and data, type `?swirl`.

3 Normalization using robust local regression

The purpose of normalization is to identify and remove sources of systematic variation, other than differential expression, in the measured fluorescence intensities (e.g. different labeling efficiencies and scanning properties of the Cy3 and Cy5 dyes; different scanning parameters, such as PMT settings; print-tip, spatial, or plate effects). It is necessary to normalize the fluorescence intensities before any analysis which involves comparing expression levels within or between slides (e.g. classification, multiple testing), in order to ensure that differences in intensities are indeed due to differential expression and not experimental artifacts. The need for normalization can be seen most clearly in self-self experiments, in which two identical mRNA samples are labeled with different dyes and hybridized to the same slide [Dudoit et al.(2002)Dudoit, Yang, Callow, and Speed]. Although there is no differential expression and one expects the red and green intensities to be equal, the red intensities often tend to be lower than the green intensities. Furthermore, the imbalance in the red and green intensities is usually not constant across the spots within and between arrays, and can vary according to overall spot intensity, location on the array, plate origin, and possibly other variables.

Location normalization. We have developed location normalization methods which correct for intensity, spatial, and other dye biases using *robust locally weighted regression* [Cleveland(1979), Yang et al.(2001)Yang, Dudoit, Luu, and Speed, Yang et al.(2002)Yang, Dudoit, Luu, Lin, Peng, Ngai, and Speed]. Local regression is a *smoothing* method for summarizing multivariate data using general curves and surfaces. The smoothing is achieved by fitting a linear or quadratic function of the predictor variables *locally* to the data, in a fashion that is analogous to computing a moving average. In the lowess and loess procedures, polynomials are fitted locally using iterated weighted least squares. *Robust* fitting guards against deviant points distorting the smoothed points. In the context of microarray experiments, robust local regression allows us to capture the non-linear dependence of the intensity log-ratio $M = \log_2 R/G$ on the overall intensity $A = \log_2 \sqrt{RG}$, while ensuring that the computed normalization values are not driven by a small number of differentially expressed genes with extreme log-ratios. For details on the R `loess` function (`modreg` package), type `?loess`.

Scale normalization. For scale normalization, a robust estimate of scale, such as the *median absolute deviation (MAD)*, may be used [Yang et al.(2001)Yang, Dudoit, Luu, and Speed, Yang et al.(2002)Yang, Dudoit, Luu, Lin, Peng, Ngai, and Speed]. For a collection of numbers x_1, \dots, x_n , the MAD is the median of their absolute deviations from the median $m = \text{median}\{x_1, \dots, x_n\}$

$$MAD = \text{median}\{|x_1 - m|, \dots, |x_n - m|\}.$$

The R function for MAD is `mad`.

Location and scale normalized intensity log-ratios M are given by

$$M \leftarrow \frac{M - l}{s},$$

where l and s denote the location and scale normalization values, respectively. The location value l can be obtained, for example, by robust local regression of M on A within print-tip-group. The scale value s could be the MAD, within print-tip-group, of location normalized log-ratios.

4 Normalization functions

4.1 Simple normalization function `maNorm`

A simple wrapper function `maNorm` is provided for users interested in applying a standard set of normalization procedures using default parameters. This function returns an object of class `marrayNorm` and has seven arguments

`mbatch`: Object of class `marrayRaw`, containing intensity data for the batch of arrays to be normalized. An object of class `marray` may also be passed if normalization is performed in several steps.

`norm`: Character string specifying the normalization procedure. Six normalization procedures are available with this function: **`none`**, for no normalization; **`median`**, for global median location normalization; **`loess`** for global intensity or A -dependent location normalization using the `loess` function; **`twoD`**, for 2D spatial location normalization using the `loess` function; **`printTipLoess`**, for within-print-tip-group intensity dependent location normalization using the `loess` function; and **`scalePrintTipMAD`**, for within-print-tip-group intensity dependent location normalization followed by within-print-tip-group scale normalization using the median absolute deviation. This argument can be specified using the first letter of each method.

`subset`: A logical or numeric vector indicating the subset of points used to compute the normalization values.

`span`: The argument `span` which controls the degree of smoothing in the `loess` function. Only used for `loess`, `twoD`, `printTipLoess`, and `scalePrintTipMAD` options.

`Mloc`: If `TRUE`, the location normalization values are stored in the slot `maMloc` of the object of class `marray` returned by the function, if `FALSE`, these values are not retained. This option allows to save memory for large datasets.

`Mscale`: If `TRUE`, the scale normalization values are stored in the slot `maMscale` of the object of class `marray` returned by the function, if `FALSE`, these values are not retained.

`echo`: If `TRUE`, the index of the array currently being normalized is printed.

4.2 Simple scale normalization function `maNormScale`

A simple wrapper function `maNormScale` is provided for users interested in applying a standard set of scale normalization procedures using default parameters. This function returns an object of class `marrayNorm` has six arguments

`mbatch`: Object of class `marrayRaw`, containing intensity data for the batch of arrays to be normalized. An object of class `marray` may also be passed if normalization is performed in several steps.

`norm`: Character string specifying the normalization procedure. Two normalization procedures are currently available for this function: `globalMAD` for global scale normalization using the median absolute deviation; `printTipMAD` for within-print-tip-group scale normalization using the median absolute deviation. This argument can be specified using the first letter of each method.

`subset`: A logical or numeric vector indicating the subset of points used to compute the normalization values.

`geo`: If `TRUE`, the MAD of each group is divided by the geometric mean of the MADs across groups [Yang et al.(2002)Yang, Dudoit, Luu, Lin, Peng, Ngai, and Speed]. This allows observations to retain their original units.

`Mscale`: If `TRUE`, the scale normalization values are stored in the slot `maMscale` of the object of class `marray` returned by the function, if `FALSE`, these values are not retained.

`echo`: If `TRUE`, the index of the array currently being normalized is printed.

The `globalMad` option, with `geo=TRUE`, allows between slide scale normalization.

4.3 General normalization function `maNormMain`

Note: We recommend users using `maNorm` and `maNormScale` rather than this function for performing standard set of normalization procedures.

This is the main internal function for location and scale normalization of cDNA microarray data is `maNormMain`; it has eight arguments (see also ? `maNormMain`):

`mbatch`: Object of class `marrayRaw`, containing intensity data for the batch of arrays to be normalized. An object of class `marrayNorm` may also be passed if normalization is performed in several steps.

`f.loc`: A list of location normalization functions, e.g., `maNormLoess`, `maNormMed`, or `maNorm2D`.

`f.scale`: A list of scale normalization functions, e.g, `maNormMAD`.

`a.loc`: For composite normalization, a function for computing the weights used in combining several location normalization functions, e.g., `maCompNormA`.

`a.scale`: For composite normalization, a function for computing the weights used in combining several scale normalization functions.

Mloc: If TRUE, the location normalization values are stored in the slot `maMloc` of the object of class `marray` returned by the function, if FALSE, these values are not retained. This option allows to save memory for large datasets.

Mscale: If TRUE, the scale normalization values are stored in the slot `maMscale` of the object of class `marray` returned by the function, if FALSE, these values are not retained.

echo: If TRUE, the index of the array currently being normalized is printed.

Normalization is performed simultaneously for each array in the batch using the location and scale normalization procedures specified by the lists of functions `f.loc` and `f.scale`. Typically, only one function is given in each list, otherwise composite normalization is performed using the weights given by `a.loc` and `a.scale` [Yang et al.(2002)Yang, Dudoit, Luu, Lin, Peng, Ngai, and Speed]. The `maNormMain` function returns objects of class `marrayNorm`.

The `marray` package contains internal functions for median (`maNormMed`), intensity or A -dependent (`maNormLoess`), and 2D spatial (`maNorm2D`) location normalization. The R robust local regression function `loess` is used for intensity dependent and 2D spatial normalization. The package also contains a function for scale normalization using the median absolute deviation (MAD) (`maNormMAD`). These functions have arguments for specifying which spots to use in the normalization and for controlling the local regression, when applicable. The functions allow normalization to be done separately within values of a layout parameter, such as `plate` or `print-tip-group`, and using different subsets of probe sequences (e.g. dilution series of control probe sequences).

5 Normalization of Swirl zebrafish microarray data

To read in the data for the Swirl experiment and generate the plate IDs

```
> library("marray", verbose = FALSE)
> data(swirl)
> maPlate(swirl) <- maCompPlate(swirl, n = 384)
```

The pre-normalization MA -plot for the Swirl 93 array in Figure 3 illustrates the non-linear dependence of the log-ratio M on the overall spot intensity A and the existence of spatial dye biases. Only a small proportion of the spots are expected to vary in intensity between the two channels. We thus perform within-print-tip-group loess location normalization using all 8,448 probes on the array.

5.1 Using simple function `maNorm`

The following command normalizes all four arrays in the Swirl experiment simultaneously. The simple wrapper function could be used to perform most of the standard normalizations procedures.

```
> swirl.norm <- maNorm(swirl, norm = "p")
> summary(swirl.norm)
```

Normalized intensity data: Object of class marrayNorm.

Call to normalization function:

```
maNormMain(mbatch = mbatch, f.loc = list(maNormLoess(x = "maA",
  y = "maM", z = "maPrintTip", w = NULL, subset = subset, span = span,
  ...)), Mloc = Mloc, Mscale = Mscale, echo = echo)
```

Number of arrays: 4 arrays.

A) Layout of spots on the array:

B) Samples hybridized to the array:

C) Summary statistics for log-ratio distribution:

	Min.	1st Qu.	Median
C:/GNU/R/R-2.4.1/library/marray/swirldata/swirl.1.spot	-2.22	-0.18	-0.01
C:/GNU/R/R-2.4.1/library/marray/swirldata/swirl.2.spot	-2.84	-0.16	0.00
C:/GNU/R/R-2.4.1/library/marray/swirldata/swirl.3.spot	-1.58	-0.23	-0.01
C:/GNU/R/R-2.4.1/library/marray/swirldata/swirl.4.spot	-2.89	-0.20	0.00
	Mean	3rd Qu.	Max.
C:/GNU/R/R-2.4.1/library/marray/swirldata/swirl.1.spot	0.05	0.20	5.09
C:/GNU/R/R-2.4.1/library/marray/swirldata/swirl.2.spot	0.00	0.16	2.17
C:/GNU/R/R-2.4.1/library/marray/swirldata/swirl.3.spot	0.04	0.25	3.24
C:/GNU/R/R-2.4.1/library/marray/swirldata/swirl.4.spot	-0.01	0.18	3.25

D) Notes on intensity data:

Spot Data

For global median normalization

```
> swirl.normm <- maNorm(swirl, norm="median")
```

5.2 Using simple function maNormScale

This simple wrapper function may be used to perform scale normalization separately from location normalization. The following examples do not represent a recommended analysis but are simply used for demonstrating the software functionality. Within-print-tip-group intensity dependent normalization followed by within-print-tip-group scale normalization using the median absolute deviation, could be performed in one step by

```
> swirl.norms <- maNorm(swirl, norm="s")
```

or sequentially by

```
> swirl.norm1 <- maNorm(swirl, norm = "p")
> swirl.norm2 <- maNormScale(swirl.norm1, norm = "p")
```

For between slide scale normalization using MAD scaled by the geometric mean of MAD across slides [Yang et al.(2001)Yang, Dudoit, Luu, and Speed, Yang et al.(2002)Yang, Dudoit, Luu, Lin, Peng, Ngai, and Spee

```
swirl.normg <- maNormScale(swirl.norm, norm="g")
```

5.3 Using main function maNormMain

The following command normalizes all four arrays in the Swirl experiment simultaneously

```
> swirl.norm <- maNormMain(swirl,  
f.loc = list(maNormLoess(x = "maA", y = "maM", z = "maPrintTip",  
w = NULL, subset = TRUE, span = 0.4)),  
f.scale = NULL,  
a.loc = maCompNormEq(),  
a.scale = maCompNormEq(),  
Mloc = TRUE, Mscale = TRUE, echo = FALSE)
```

This is the default normalization procedure in `maNormMain`, thus the same results could be obtained by calling

```
> swirl.norm <- maNormMain(swirl)
```

To see the effect of within-print-tip-group location normalization, compare the pre-and post-normalization boxplots and MA -plots in Figures 1, 2, and 3. Normalized log-ratios M are now evenly distributed about about zero across the range of intensities A for each print-tip-group. Furthermore, the non-linear location normalization seems to have eliminated, to some extent, the scale differences among print-tip-groups and arrays.

5.4 Plots

The plots were produced using the following commands:

```
> boxplot(swirl[, 3], xvar = "maPrintTip", yvar = "maM", main = "Swirl array 93: pre--normaliza  
> boxplot(swirl, yvar = "maM", main = "Swirl arrays: pre--normalization")  
> boxplot(swirl.norm[, 3], xvar = "maPrintTip", yvar = "maM", main = "Swirl array 93: post--no  
> boxplot(swirl.norm, yvar = "maM", main = "Swirl arrays: post--normalization")  
> plot(swirl[, 3], main = "Swirl array 93: pre--normalization MA--plot")  
> plot(swirl.norm[, 3], main = "Swirl array 93: post--normalization MA--plot")
```

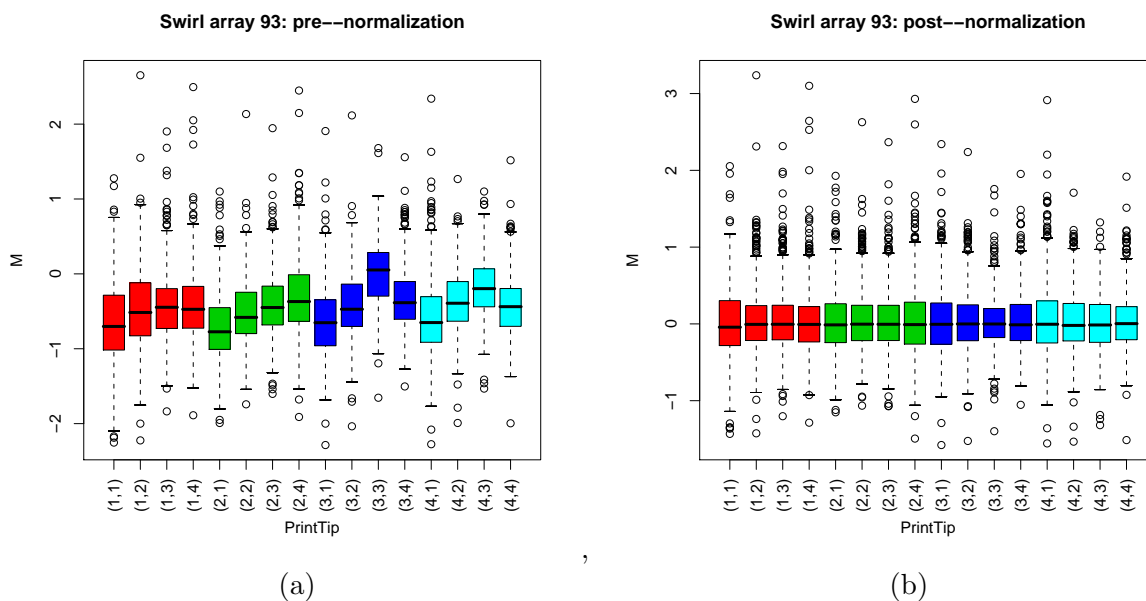


Figure 1: Boxplots by print-tip-group of the pre- and post-normalization intensity log-ratios M for the Swirl 93 array.

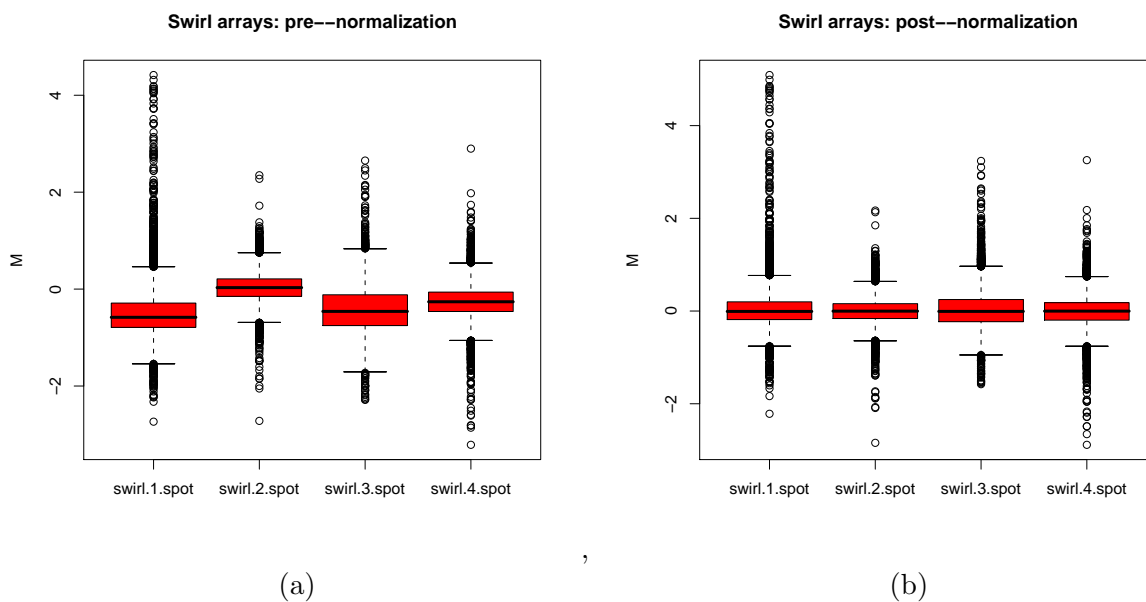


Figure 2: Boxplots of the pre- and post-normalization intensity log-ratios M for the four arrays in the Swirl experiment.

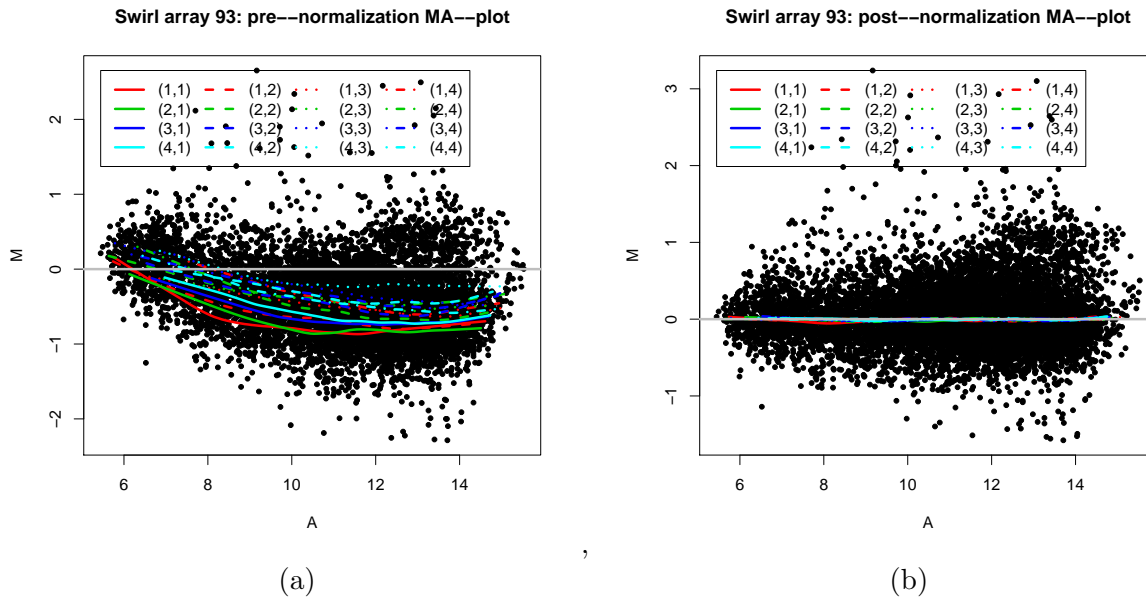


Figure 3: Pre- and post-normalization MA -plot for the Swirl 93 array, with the lowest fits for individual print-tip-groups. Different colors are used to represent lowest curves for print-tips from different rows, and different line types are used to represent lowest curves for print-tips from different columns.

Note: Sweave. This document was generated using the `Sweave` function from the `R tools` package. The source file is in the `/inst/doc` directory of the package `marray`.

References

- [Cleveland(1979)] W. S. Cleveland. Robust locally weighted regression and smoothing scatterplots. *Journal of the American Statistical Association*, 74(368):829–836, 1979.
- [Dudoit and Yang(2002)] S. Dudoit and Y. H. Yang. Bioconductor R packages for exploratory analysis and normalization of cDNA microarray data. In G. Parmigiani, E. S. Garrett, R. A. Irizarry, and S. L. Zeger, editors, *The Analysis of Gene Expression Data: Methods and Software*. Springer, New York, 2002.
- [Dudoit et al.(2002)Dudoit, Yang, Callow, and Speed] S. Dudoit, Y. H. Yang, M. J. Callow, and T. P. Speed. Statistical methods for identifying differentially expressed genes in replicated cDNA microarray experiments. *Statistica Sinica*, 12(1):111–139, 2002.
- [Yang et al.(2001)Yang, Dudoit, Luu, and Speed] Y. H. Yang, S. Dudoit, P. Luu, and T. P. Speed. Normalization for cDNA microarray data. In M. L. Bittner, Y. Chen, A. N. Dorsel, and E. R. Dougherty, editors, *Microarrays: Optical Technologies and Informatics*, volume 4266 of *Proceedings of SPIE*, May 2001.

[Yang et al.(2002)Yang, Dudoit, Luu, Lin, Peng, Ngai, and Speed] Y. H. Yang, S. Dudoit, P. Luu, D. M. Lin, V. Peng, J. Ngai, and T. P. Speed. Normalization for cDNA microarray data: a robust composite method addressing single and multiple slide systematic variation. *Nucleic Acids Research*, 30(4), 2002.