

Bioconductor's marrayPlots package

Sandrine Dudoit¹ and Yee Hwa Yang²

November 25, 2003

1. Division of Biostatistics, University of California, Berkeley,

<http://www.stat.berkeley.edu/~sandrine>

2. Department of Statistics, University of California, Berkeley, yeehwa@stat.berkeley.edu

Contents

1 Overview	1
2 Getting started	2
3 Diagnostic plots	2
4 Spatial plots of spot statistics – maImage	3
5 Boxplots of spot statistics – maBoxplot	8
6 Scatter-plots of spot statistics – maPlot	11
7 Wrapper functions for basic sets of diagnostic plots	13

1 Overview

This document provides a detailed discussion of the **marrayPlots** package, which is part of a suite of four packages for diagnostic plots and normalization of cDNA microarray data. The package is given in **marrayPlotsShort**. Information on the other packages can be found in Dudoit and Yang (2002) and the vignettes for each package. Like most Bioconductor packages, these four packages rely on the object-oriented class/method mechanism, provided by the R **methods** package, to allow efficient and systematic representation and manipulation of microarray data. **marrayPlots** provides functions for diagnostic plots of microarray spot statistics, such as boxplots, scatter-plots, and spatial color images. Examination of diagnostic plots of intensity data is important in order to identify printing, hybridization, and scanning artifacts which can lead to biased inference concerning gene expression. The other three packages are

marrayClasses. This package contains basic class definitions and associated methods for pre- and post-normalization intensity data for batches of arrays.

marrayInput. This package provides functionality for reading microarray data into R, such as intensity data from image processing output files (e.g. **.spot** and **.gpr** files for the **Spot** and

GenePix packages, respectively) and textual information on probes and targets (e.g. from gal files and god lists). **tcltk** widgets are supplied to facilitate and automate data input and the creation of microarray specific R objects for storing these data.

marrayNorm. 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

Installing the package. To install the **marrayPlots** package for Windows operating systems, first download the file **marrayPlots-snapshot.zip** from the Bioconductor website <http://www.bioconductor.org/packages/html/marrayPlots.html>. Next, after starting R, from the menu select **Packages**, then **Install package from local zip file...** Find and highlight the location of the zip file and click on **open**.

Loading the package. To load the **marrayPlots** package in your R session, type `library(marrayPlots)`.

Help files. As with any R package, detailed information on functions, classes and methods can be obtained in the help files. For instance, to view the help file for the function **maImage** in a browser, use `help.start()` followed by `? maImage`.

Microarray classes. The **marrayPlots** packages relies on microarray class definitions in **marrayClasses**. You should also install this package and consult its vignette for more information.

Case study. We demonstrate the functionality of this collection of R packages using gene expression data from the Swirl zebrafish experiment. These data are included as part of the **marrayInput** 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`.

Demo. Code for a demo of the package is in the `/demo` directory. To run the demo, type `demo(marrayPlots)`.

Next. After reading your data into R using **marrayInput** and producing diagnostic plots using **marrayPlots**, the **marrayNorm** package can be used for normalization of the fluorescence intensities.

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 **marrayPlots**.

3 Diagnostic plots

Before proceeding to normalization or any higher-level analysis, it is instructive to look at diagnostic plots of spot statistics, such as red and green foreground and background log-intensities, intensity

log-ratio, area, etc. Such plots are useful for the purpose of identifying printing, hybridization, and scanning artifacts as demonstrated below. Three main types of functions were defined to operate on pre- and post-normalization microarray objects: functions for boxplots, scatter-plots, and spatial images. The main arguments to these functions are microarray objects of classes `marrayRaw`, `marrayNorm`, or `marrayTwo`, and arguments specifying which spot statistics to display (e.g. Cy3 and Cy5 background intensities, intensity log-ratios) and which subset of spots to include in the plots. Default graphical parameters are chosen for convenience using the function `maDefaultPar` (e.g. color palette, axis labels, plot title), but the user has the option to overwrite these parameters at any point. Note that by default the plots are done for the first array in a batch. To produce plots for other arrays, subsetting methods may be used. For example, to produce diagnostic plots for the second array in the batch of zebrafish arrays `swirl`, the argument `swirl[,2]` should be passed to the plot functions.

To read in the data for the Swirl experiment and generate the plate IDs (see `marrayClasses` and `marrayInput` for greater details)

```
> library(marrayPlots, verbose = FALSE)
```

Welcome to Bioconductor

Vignettes contain introductory material. To view,
simply type: `openVignette()`
For details on reading vignettes, see
the `openVignette` help page.

```
> library(Biobase, verbose = FALSE)
```

```
> library(marrayNorm, verbose = FALSE)
```

Loading required package: stepfun

```
> data(swirl)
```

```
> maPlate(swirl) <- maCompPlate(swirl, n = 384)
```

4 Spatial plots of spot statistics – `maImage`

The function `maImage` creates *images* of shades of gray or colors that correspond to the values of a statistic for each spot on an array. Details on the arguments of the function are given in `?maImage`. The statistic can be the intensity log-ratio M , a spot quality measure (e.g. spot size or shape), or a test statistic. This function can be used to explore whether there are any spatial effects in the data, for example, print-tip or cover-slip effects. In addition to existing color palette functions, such as `rainbow` and `heat.colors`, a new function `maPalette` was defined to generate color palettes from user supplied low, middle, and high color values. To create white-to-green, white-to-red, and green-to-red palettes for microarray images

```
> Gcol <- maPalette(low = "white", high = "green", k = 50)
```

```
> Rcol <- maPalette(low = "white", high = "red", k = 50)
```

```
> RGcol <- maPalette(low = "green", high = "red", k = 50)
```

Useful diagnostic plots are images of the Cy3 and Cy5 background intensities; these images may reveal hybridization artifacts such as scratches on the slides, drops, cover-slip effects etc. The following commands produce images of the Cy3 and Cy5 background intensities for the Swirl 93 array (third array in the batch) using white-to-green and white-to-red color palettes, respectively.

```
> tmp <- maImage(swirl[, 3], x = "maGb", subset = TRUE, col = Gcol,
+   contours = FALSE, bar = FALSE)

> tmp <- maImage(swirl[, 3], x = "maRb", subset = TRUE, col = Rcol,
+   contours = FALSE, bar = FALSE)
```

Note that the same images can be obtained using the default arguments of the function by the shorter commands

```
maImage(swirl[,3], x="maGb")
maImage(swirl[,3], x="maRb")
```

If `bar=TRUE`, a calibration color bar is displayed to the right of the images. The `maImage` function returns the values and corresponding colors used to produce the color bar, as well as a six number summary of the spot statistics. The resulting images are shown in Figure 1. It can be noted that the Cy3 and Cy5 background intensities are not uniform across the slide and are higher in the top right corner, perhaps due to cover slip effects or tilt of the slide during scanning. Such patterns were not as clearly visible in the individual Cy3 and Cy5 TIFF images. Similar displays of the Cy3 and Cy5 foreground intensities do not exhibit such strong spatial patterns. For other arrays, such as the Swirl 81 array, background images revealed the existence of a scratch with very high background in print-tip-groups (3,2) and (3,3).

The `maImage` function may also be used to generate an image of the pre-normalization log-ratios M (or any other statistic of interest), using a green-to-red color palette. Figure 2 displays such an image for the Swirl 93 array, highlighting only those spots with the highest and lowest 10% pre-normalization log-ratios M . Other options include displaying contours and altering graphical parameters such as axis labels and plot title. Figure 2 suggests the existence of spatial dye biases in the intensity log-ratio, with higher values in grid (3,3) and lower values in grid column 1 of the array.

```
> tmp <- maImage(swirl[, 3], x = "maM", bar = FALSE, main = "Swirl array 93: image of pre--norm

> tmp <- maImage(swirl[, 3], x = "maM", subset = maTop(maM(swirl[,
+   3])), h = 0.1, l = 0.1), col = RGcol, contours = FALSE, bar = FALSE,
+   main = "Swirl array 93: image of pre--normalization M for % 10 tails")
```

Note that the `maImage` function (and other functions `maBoxplot` and `maPlot` to be described next) can be used to plot other statistics than fluorescence intensities. They can be used to plot layout parameters such as spot coordinates `maSpotRow`, print-tip-group coordinates `maPrintTip`, or plate IDs `maPlate` (Figure 3).

```
> tmp <- maImage(swirl[, 3], x = "maSpotCol", bar = FALSE)

> tmp <- maImage(swirl[, 3], x = "maPrintTip", bar = FALSE)
```

```
> tmp <- maImage(swirl[, 3], x = "maControls", col = heat.colors(10),  
+   bar = FALSE)  
  
> tmp <- maImage(swirl[, 3], x = "maPlate", bar = FALSE)
```

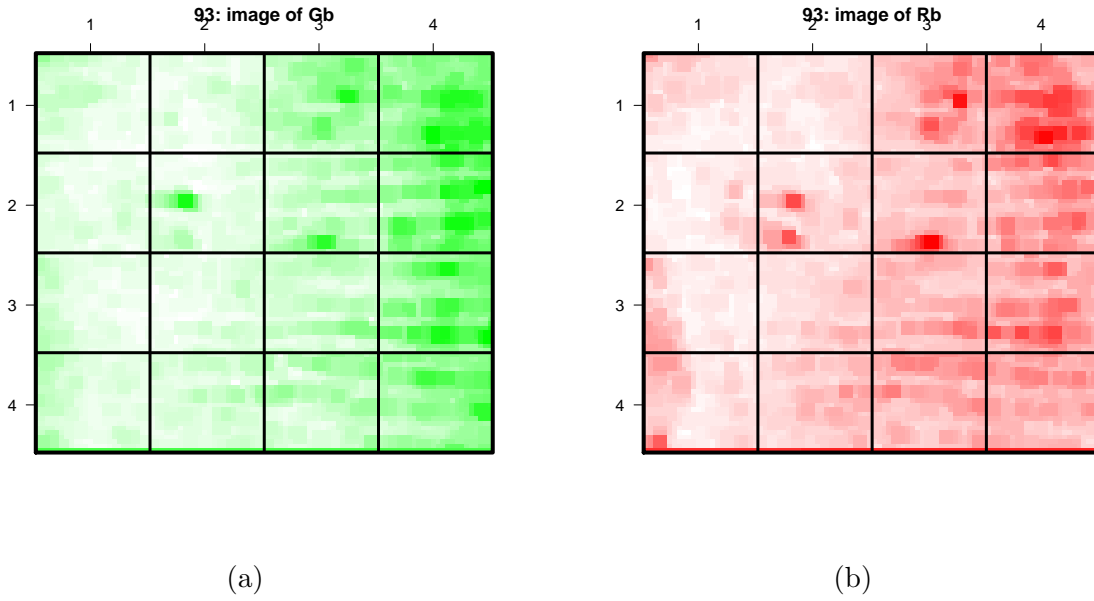


Figure 1: Images of background intensities for the Swirl 93 array. Panel (a): Cy3 background intensities using white-to-green color palette. Panel (b): Cy5 background intensities using white-to-red color palette.

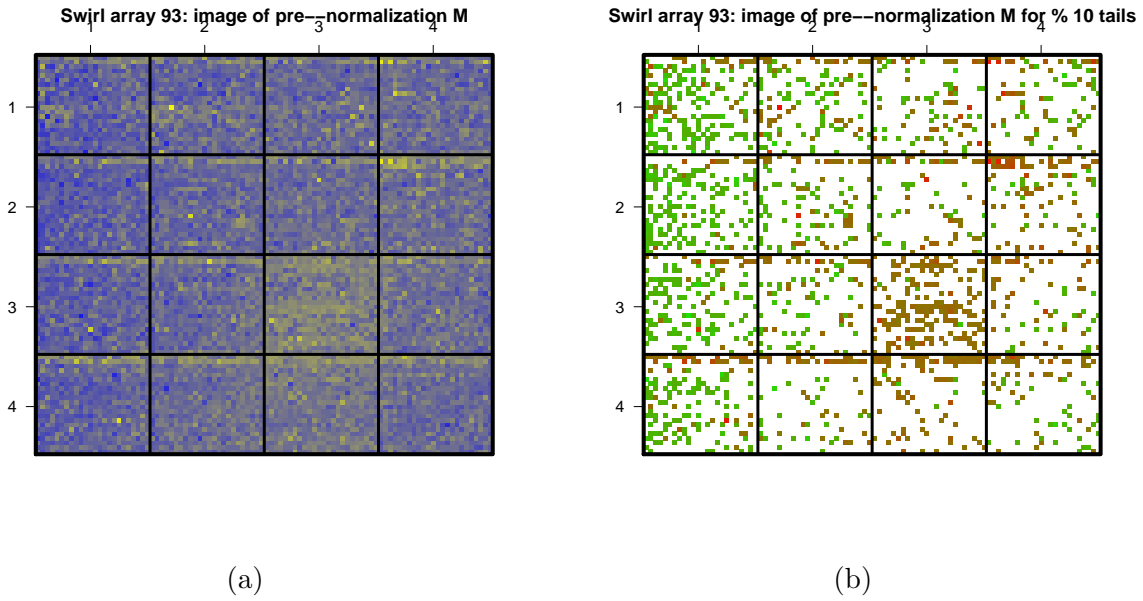


Figure 2: Images of the pre-normalization intensity log-ratios M for the Swirl 93 array, using a green-to-red color palette. Panel (a): All spots are displayed. Panel (b): only spots with the highest and lowest 10% log-ratios are highlighted.

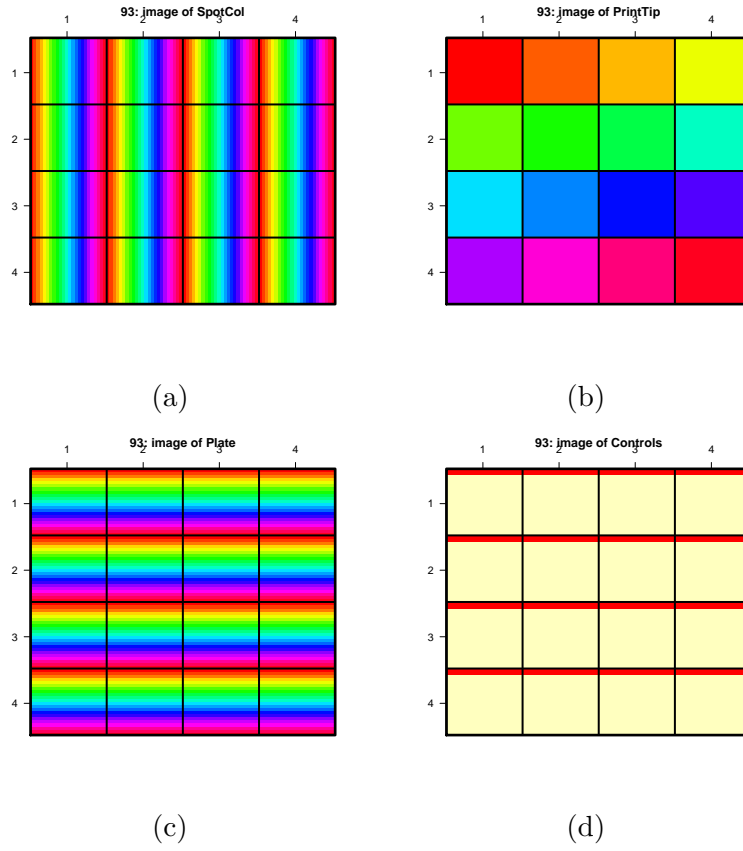


Figure 3: Images of layout parameters for the Swirl 93 array. Panel (a): Spot matrix column coordinate. Panel (b): Print-tip-group. Panel (c): Plate index. Panel (d): Control status.

5 Boxplots of spot statistics – maBoxplot

Boxplots of spot statistics by plate, print-tip-group, or slide can also be useful to identify spot or hybridization artifacts. *Boxplots*, also called *box-and-whisker plots*, were first proposed by Tukey in 1977 as simple graphical summaries of the distribution of a variable. The summary consists of the median, the upper and lower quartiles, the range, and, possibly, individual extreme values. The central box in the plot represents the *inter-quartile range (IQR)*, which is defined as the difference between the 75th percentile and 25th percentile, i.e., the upper and lower quartiles. The line in the middle of the box represents the median; a measure of central location of the data. Extreme values, greater than 1.5 IQR above the 75th percentile and less than 1.5 IQR below the 25th percentile, are typically plotted individually.

The function `maBoxplot` produces boxplots of microarray spot statistics for the classes `marrayRaw`, `marrayNorm`, and `marrayTwo` (see details in `? maBoxplot`). The function `maBoxplot` has three main arguments

- `m`: Microarray object of class `marrayRaw`, `marrayNorm`, or `marrayTwo`.
- `x`: Name of accessor method for the spot statistic used to stratify the data, typically a slot name for the microarray layout object such as `maPlate` or a method such as `maPrintTip`. If `x` is `NULL`, the data are not stratified.
- `y`: Name of accessor method for the spot statistic of interest, typically a slot name for the microarray object `m`, such as `maM`.

Figure 4 panel (a) displays boxplots of pre-normalization log-ratios M for each of the 16 print-tip-groups for the Swirl 93 array. This plot was generated by the following commands

```
> maBoxplot(swirl[, 3], x = "maPrintTip", y = "maM", main = "Swirl array 93: pre--normalization")
```

The boxplots clearly reveal the need for normalization, since most log-ratios are negative in spite of the fact that only a small proportion of genes are expected to be differentially expressed in the mutant and wild-type zebrafish. As is often the case, this corresponds to higher signal in the Cy3 channel than in the Cy5 channel even in the absence of differential expression. In addition, the boxplots show the existence of spatial dye biases in the log-ratios. In particular, print-tip-group (3,3) clearly stands out from the remaining ones, as suggested also in the image of Figure 2. The function `maBoxplot` may also be used to produce boxplots of spot statistics for all arrays in a batch. Such plots are useful when assessing the need for between array normalization, for example, to deal with scale differences among different arrays. The following command produces a boxplot of the pre-normalization intensity log-ratios M for each array in the batch `swirl`. Figure 5 panel (a) suggest that different normalizations may be required for different arrays, including possibly scale normalization.

```
> maBoxplot(swirl, y = "maM", main = "Swirl arrays: pre--normalization")
```

The function `maNorm` from the `marrayNorm` package can be used for within-print-tip-group loess location normalization using all 8,448 probes on the array. The following command normalizes all four arrays in the Swirl experiment simultaneously. Please refer to the vignette or help files of the `marrayNorm` package for more information (e.g. `? maNorm`).

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

The following commands can be used to produce post-normalization boxplots of the log-ratios. The plots are shown in panel (b) of Figures 4 and 5.

```
> maBoxplot(swirl.norm[, 3], x = "maPrintTip", y = "maM", main = "Swirl array 93: post--normalization")
```

```
> maBoxplot(swirl.norm, y = "maM", col = "green", main = "Swirl arrays: post--normalization")
```

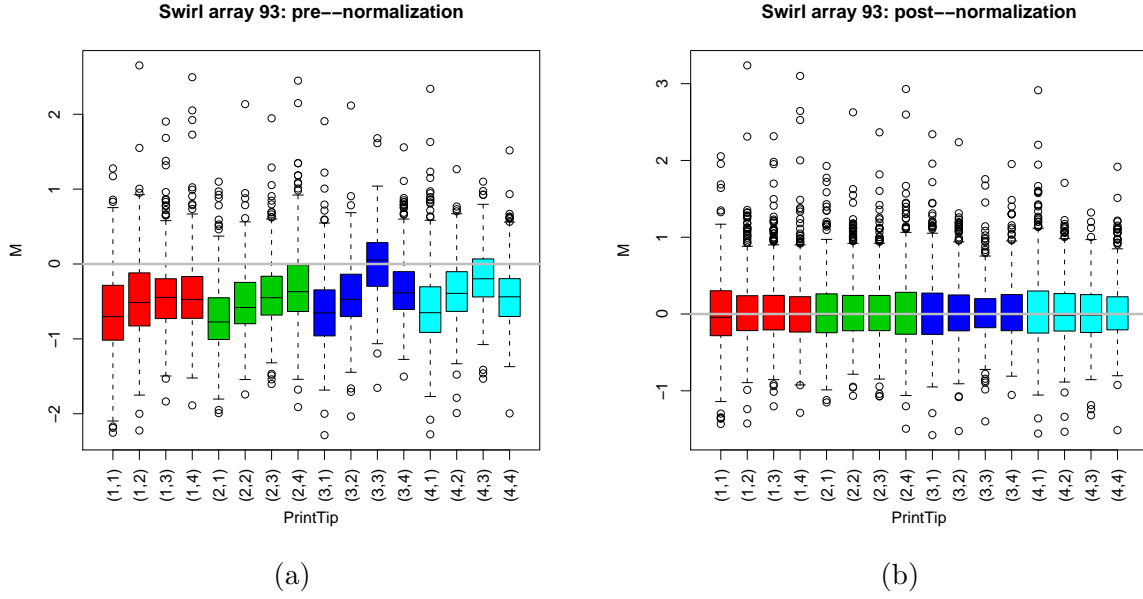


Figure 4: Boxplots by print-tip-group of the pre- and post-normalization intensity log-ratios M for the Swirl 93 array.

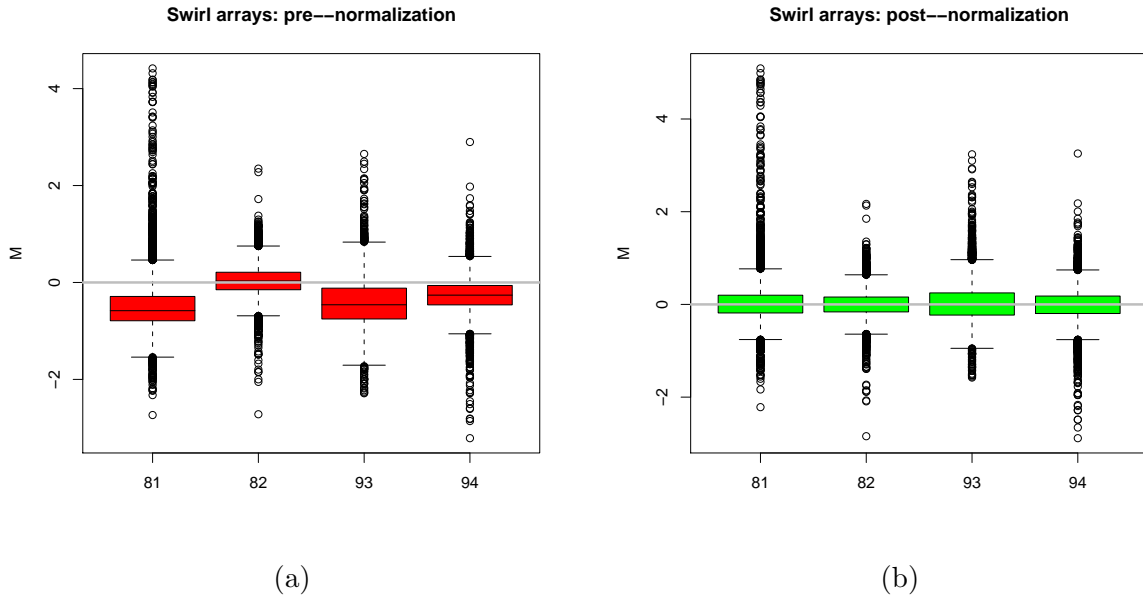


Figure 5: Boxplots of the pre- and post-normalization intensity log-ratios M for the four arrays in the Swirl experiment.

6 Scatter-plots of spot statistics – maPlot

The function `maPlot` produces *scatter-plots* of microarray spot statistics for the classes `marrayRaw`, `marrayNorm`, and `marrayTwo`. It also allows the user to highlight and annotate subsets of points on the plot, and display fitted curves from robust local regression or other smoothing procedures (see details in ? `maPlot`). The function `maPlot` has seven main arguments

- `m`: Microarray object of class `marrayRaw`, `marrayNorm`, or `marrayTwo`.
- `x`: Name of accessor function for the abscissa spot statistic, typically a slot name for the microarray object `m`, such as `maA`.
- `y`: Name of accessor function for the ordinate spot statistic, typically a slot name for the microarray object `m`, such as `maM`.
- `z`: Name of accessor method for the spot statistic used to stratify the data, typically a slot name for the microarray layout object such as `maPlate` or a method such as `maPrintTip`. If `z` is `NULL`, the data are not stratified.
- `lines.func`: Function for computing and plotting smoothed fits of `y` as a function of `x`, separately within values of `z`, e.g. `maLoessLines`. If `lines.func` is `NULL`, no fitting is performed.
- `text.func`: Function for highlighting a subset of points, e.g., `maText`. If `text.func` is `NULL`, no points are highlighted.
- `legend.func`: Function for adding a legend to the plot, e.g. `maLegendLines`. If `legend.func` is `NULL`, there is no legend.

As usual, optional graphical parameters may be supplied and these will overwrite the default parameters set in the plot functions. A number of functions for computing and plotting the fits are provided in `marrayPlot`, such as `maLowessLines` and `maLoessLines` for robust local regression using the R functions `lowess` and `loess`, respectively (type ? `loess` or ? `lowess` for a brief description of R functions for robust local regression). Functions are also provided for highlighting points (e.g. `maText`) and adding a legend to the plot (e.g. `maLegendLines`).

MA-plots. Single-slide expression data are typically displayed by plotting the log-intensity $\log_2 R$ in the red channel vs. the log-intensity $\log_2 G$ in the green channel. Such plots tend to give an unrealistic sense of concordance between the red and green intensities and can mask interesting features of the data. We thus recommend plotting the intensity log-ratio $M = \log_2 R/G$ vs. the mean log-intensity $A = \log_2 \sqrt{RG}$. An *MA*-plot amounts to a 45° counterclockwise rotation of the $(\log_2 G, \log_2 R)$ -coordinate system, followed by scaling of the coordinates. It is thus another representation of the (R, G) data in terms of the log-ratios M which directly measure differences between the red and green channels and are the quantities of interest to most investigators. We have found *MA*-plots to be more revealing than their $\log_2 R$ vs. $\log_2 G$ counterparts in terms of identifying spot artifacts and for normalization purposes (Dudoit et al., 2002; Yang et al., 2001, 2002).

Figure 6 panel (a) displays the pre-normalization *MA*-plots for the Swirl 93 array, with the sixteen lowess fits for each of the print-tip-groups (using a smoother span $f = 0.3$ for the `lowess` function). The figure was generated with the following commands

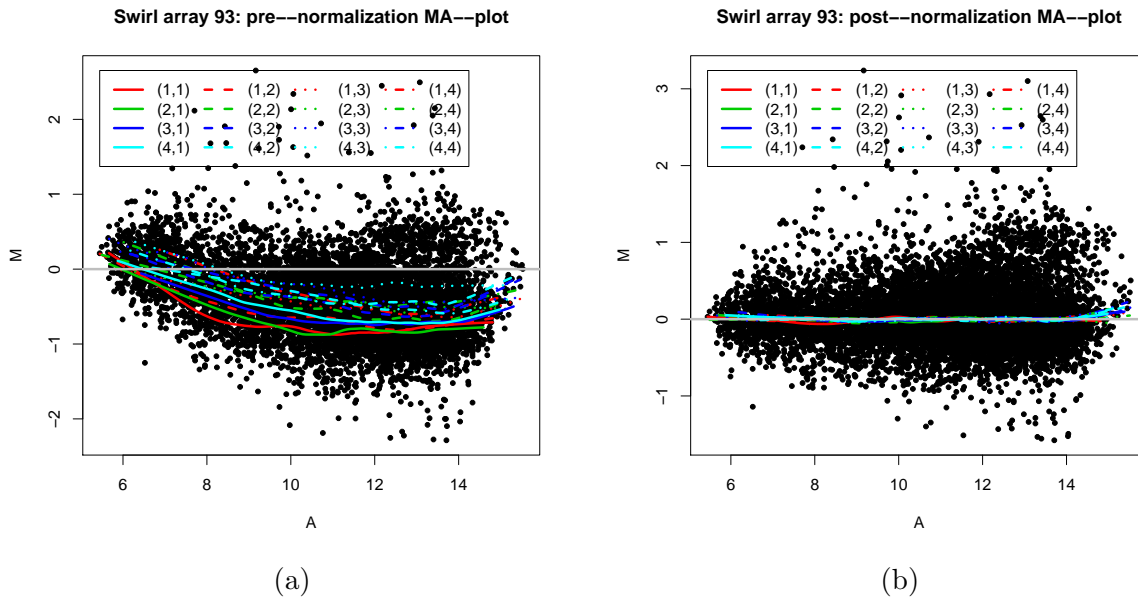


Figure 6: Pre- and post-normalization MA -plot for the Swirl 93 array, with the lowess fits for individual print-tip-groups. Different colors are used to represent lowess curves for print-tips from different rows, and different line types are used to represent lowess curves for print-tips from different columns.

```
> defs <- maDefaultPar(swirl[, 3], x = "maA", y = "maM", z = "maPrintTip")
> legend.func <- do.call("maLegendLines", defs$def.legend)
> lines.func <- do.call("maLowessLines", c(list(TRUE, f = 0.3),
+   defs$def.lines))
> maPlot(swirl[, 3], x = "maA", y = "maM", z = "maPrintTip", lines.func,
+   text.func = maText(), legend.func, main = "Swirl array 93: pre--normalization MA--plot")

> maPlot(swirl.norm[, 3], x = "maA", y = "maM", z = "maPrintTip",
+   lines.func, text.func = maText(), legend.func, main = "Swirl array 93: post--normalization MA--plot")
```

The same plots can be obtain using the default arguments of the function by the commands

```
maPlot(swirl[,3])
maPlot(swirl.norm[,3])
```

To highlight, say, the spots with the highest and lowest 5% log-ratios using purple symbols "0", set `text.func=maText(subset=maTop(maM(swirl[,3]),h=0.05,l=0.05),labels="0",col="purple")`. Figure 6 illustrates the non-linear dependence of the log-ratio M on the overall spot intensity A and thus suggests that an intensity or A -dependent normalization method is preferable to a global one (e.g. median normalization). Also, the lowess fits vary among print-tip-groups, again revealing the existence of spatial dye biases. Figure 6 panel (b) displays the MA -plot after within-print-tip-group loess location normalization.

7 Wrapper functions for basic sets of diagnostic plots

Three wrapper functions are provided to automatically generate a standard set of diagnostic plots: functions `maDiagnPlots1`, `maRawPlots`, and `maNormPlots`. For example, `maDiagnPlots1` produces eight plots of pre- and post-normalization cDNA microarray data: color images of Cy3 and Cy5 background intensities, and of pre- and post-normalization log-ratios M ; boxplots of pre- and post-normalization log-ratios M by print-tip-group; MA -plots of pre- and post-normalization log-ratios M by print-tip-group. All three functions provide options for saving the figures to a file, in postscript or jpeg format.

```
maDiagnPlots1(swirl[,2], title="Swirl array 93: Diagnostic plots",
save=TRUE, fname="swirl93.jpeg", dev="jpeg")
```

References

- 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.
- 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.
- 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.
- 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.