



*Proceedings of the 3rd International Workshop
on Distributed Statistical Computing (DSC 2003)
March 20–22, Vienna, Austria ISSN 1609-395X
Kurt Hornik, Friedrich Leisch & Achim Zeileis (eds.)
<http://www.ci.tuwien.ac.at/Conferences/DSC-2003/>*

A Graphical Users Interface to Normalize Microarray Data

Fatima Sanchez Cabo^{‡§}, Zlatko Trajanoski[§], Kwang-Hyun Cho*,
Olaf Wolkenhauer[†]

[‡]Department of Biomolecular Sciences, UMIST, Manchester, U.K.

[§]Institute of Biomedical Engineering, Graz University of Technology, Graz, Austria

* School of Electrical Engineering, University of Ulsan, Ulsan, South Korea

[†]Department of Computer Science, University of Rostock, Rostock, Germany

Abstract

Microarray technology is becoming an essential tool in functional genomics. The possibility of monitoring the expression level of thousands of genes simultaneously, as the response to a particular biological condition, gives to the biologists the chance to widen the aims of their experiments and opens a door to the understanding of cellular transcription processes. In order to extract valuable information from the big amount of data that microarrays experiments generate, suitable and powerful statistical and computational methods are required. An example of the effort of statisticians and computer scientists is the release of the first Bioconductor software and the increasing number of functions for microarray data analysis implemented in several programming languages (e.g. R, MATLAB, Java) by different research teams all around the world.

In this paper, we describe a Graphical Users Interface (GUI) written in MATLAB to deal with the normalization of microarray data. In our opinion, not enough importance has been given yet either to the assessment of the effect of the normalization on the data or to the study of the most suitable normalization methods according to the experimental design. To aim these objectives, a great variety of normalization methods were implemented in the interface here described, allowing the user to visualize the data before and after every step of the normalization process. Our interface suggests an example of what should be done using also other softwares such as R.

The features implemented in this interface were validated using data sets from microarray experiments carried out for *Mycobacterium tuberculosis* by

the Bacterial Microarray Group St. George's Hospital, Medical School in London and for *Streptomyces coelicolor* by the Streptomyces group at UMIST.

1 Introduction

Two color microarrays measure the relative abundance of messenger RNA (mRNA) of thousands of genes in two different samples. To obtain an estimator of the mRNA abundance, the two pools of mRNA from the cell populations to be studied are reverse transcribed to complementary DNA (cDNA) and labelled using two different fluorescent dyes (usually cyanine dyes Cy3 and Cy5), as described in Eisen and Brown (1999) and Schulze and Downward (2001). The two pools are then combined and applied to the microarray itself, where products of the polymerase chain reaction (PCR) generated from cDNA libraries or clone collections were printed as spots at defined locations. Labelled cDNA or genomic DNA (gDNA) in the pools hybridize to complementary sequences on the array and unhybridized DNA is washed off. The slide is then scanned using two different wavelengths and the intensity of the same spot in both channels is compared. This results in a measurement of the ratio of transcript levels for each gene represented on the array.

The statistical analysis starts with the scanning file itself. Different location parameters for the distribution of the pixels in a particular spot are given (e.g., mean, median, mode) and the most suitable one to explain the intensity value of a given spot in both channels should be chosen. The scan file gives the location parameters for both channels foreground intensities (Cy3 and Cy5) and for their background. The background intensity measures the intensity of the mRNA that binds to the slide even if there is no material spotted. Using all this information the next step is to filter those spots with bad quality that should not be used for further analysis. Yet, before proper analysis of the data they need to be normalized in order to remove the non-biological variation introduced by the experimental process and to enable the comparison of the intensity values within and across slides. However, these correction methods might introduce additional noise or could even falsely transform the data if the assumptions they made are not carefully observed. The GUI here described can help in the detection of over-fitting or additional noise introduced.

There are many different methods to normalize microarray data. Some statistical and algebraic methods such as ANOVA (described in Kerr and Churchill (2001); Kerr, Martin, and Churchill (2000)) and SVD (see Alter, Brown, and Botstein (2000)) can be applied to normalize the data. These methods aim to remove the non-biological variation in one single step. However, they can be considered by biologists as "black-boxes". In consequence, it can be of a greater interest the use of a sequential method, allowing the user to choose different options at different stages of the normalization process, according to the particularities of the experiment. The interface here described normalizes the data in this way. Table 1 is an example of the very general sequential method implemented in our toolbox MADE (MicroArray Data Explorer).

The paper is organized as follows: Firstly, the motivation for normalization of microarray data is explained and the main sources of variability in microarray data are defined. These sources of variability can be introducing in the data biological variability but also random and systematic errors. The main sources of non-biological variation are the background effect, the dye effect and the array ef-

Table 1: General approach to the sequential process. All the options were implemented in our toolbox in order to enable the user to choose the most suitable one according to its experimental design.

NORMALIZATION IN MADE (MicroArray Data Explorer)		
<i>Effects corrected</i>	<i>Options</i>	
<i>Background effect</i>	<ol style="list-style-type: none"> 1. Background subtraction 2. No subtraction 	
<i>Spatial effect</i>	$p_i = \frac{r_i}{g_i}$	
<i>Dye effect</i>	<i>Using all genes</i>	<ol style="list-style-type: none"> 1. Global constant 2. Linear regression 3. LOWESS function 4. LOWESS for print-tips
	<i>Quality control elements</i>	<ol style="list-style-type: none"> 1. Dye-swap normalization 2. Use of spotted controls
<i>Array effect: Across replicates normalization</i>		
<i>Average experimental replicates (slides/spots)</i>		
<i>Array effect: Across samples normalization</i>	<ol style="list-style-type: none"> 1. Against all arrays 2. Against arrays in J 	
<i>Transformation of the data</i>	<ol style="list-style-type: none"> 1. $\log_2(\bullet)$ transformation 2. $\sqrt{\bullet}$ transformation 3. $\text{lin-}\log_2$ transformation 4. $\text{arsinh}(\bullet)$ transformation 	

fect. The paper describes then the different features that allow the user to visualize and correct every of those effects.

2 Normalization of microarray data

A proper understanding of the intrinsic errors in a measurement requires a suitable mathematical approach. Errors in a measurement can be of two types: systematic or random.

- *Systematic error or inaccuracy* is a fixed positive or negative error that is the same if the measurement is repeated (systematic error).
- *Random error or imprecision* is a random positive or negative error that varies every time the measurement is made.

In the microarray production process, many systematic and random errors are introduced. These errors are masking the biological variation in which we are interested. Normalization is the process of removing all this non-biological variation. As described in Section 1, microarrays are tools used to estimate the amount of mRNA for a gene across different conditions. For such an estimation, we rely on two intensity values per spot, one for every channel. In a very simple approach, according to [Kepler, Crosby, and Morgan \(2002\)](#), every intensity value can be modelled as:

$$I = N \cdot A + error$$

where A is the abundance of mRNA for the gene in the given sample. N is the normalization factor that corrects all systematic errors and $error$ summarizes the random error. The objective of the normalization process is to make I a reliable estimator of A . For that it is essential to estimate N and $error$. According to [Kerr and Churchill \(2001\)](#) and [Kerr et al. \(2000\)](#), there are four main sources of variation in microarray data. Some of them introduce non-biological variability (contributing to N and $error$) and should be removed in order to understand the real biological variation. The four main sources of variation are:

- Dye effect. The different incorporation properties of the dyes and their different physical characteristics make this the most important source of systematic error in two-color microarrays.
- Array effect. The difference in the overall intensity across different arrays can be due to real biological variation from one condition to another or just to some experimental noise.
- Gene effect. The different expression level of a particular gene in a particular array can be due to the biological variability of the gene or to some noise.
- Sample effect. If the overall intensity of the hybridized samples is different, it can be due to some experimental error or to real biological activity.

Besides these four factors the background effect must be also considered. Some part of the probe will attach to the slide even when there is not spotted material, contributing to the foreground intensity. Some efforts are being done to provide a reliable estimator for the background intensity, as shown in [Kooperberg, Fazio, J.J., and Tsukiyama \(2002\)](#).

The increasing number of methods described to correct all the systematic and random bias mostly summarized in the four effects previously described (see [Yang, Dudoit, Lin, Peng, Ngai, and Speed \(2002\)](#)) may lead to confusion in the analysis of microarray data. Which method should I use? Should I use all of them? These are typical questions when facing the normalization of a data set. With the aim of helping to answer these questions, we implemented an interface in MATLAB[©] (Mathworks Inc.)

3 A normalization interface implemented in MATLAB[©]

We chose MATLAB to analyze the data from our microarray experiments due, among other reasons, to the variety of representation features that this software provides. However, MATLAB has many limitations in terms of memory and speed and many important statistical tests are not implemented in its library. For example, it is limited in the functions for multivariate analysis of variance, which is becoming an increasingly important tool for the normalization and analysis of microarray data. For all those purposes, programs based on R, such as Bioconductor, or functions written in Java or C++ can be more appropriated. In this direction, the Jackson Laboratory [Churchill \(2002\)](#) has implemented a number of functions in MATLAB with C++ core functions, improving the efficiency of functions such as LOWESS by [Cleveland \(1979\)](#) and implementing factorial designs that were not written as default MATLAB functions. Some ideas about the implementation of GUI interfaces in R has been presented in [Unwin \(2001\)](#).

One of the main problems in the normalization of microarray data is the organization of all the information. The MATLAB interface described in this paper summarizes the different steps that must be performed in the normalization of microarray data, allowing the user to visualize different plots in order to decide at every stage which is the most suitable option among all those available. A compromise between particularity and generalization must be taken in the normalization process. For this reason, although the methods must be as general as possible to allow its application to all kind of data sets, it is essential to visualize the data to define particularities associated with it.

With our interface, we tried to offer a wide sample of methods to correct the effects previously defined. Among them, all the gene specific errors (e.g. short PCR products) and the spatial noise affect both channels in the same amount. Consequently, these sources of noise are mainly removed just by taking the ratio of both channels. The new version of the GUI will include a spatial normalization based in the method by [Colantuoni, Henry, Zeger, and Pevsner \(2002\)](#). Hence, the interface has three main blocks to correct the background, dye and array effects. For every of them the interface enables:

1. Visualization of different plots for the pre-corrected data in order to choose the most suitable normalization method.
2. A variety of effect-correction options to be chosen.
3. Visualization of the corrected data in order to asses the effect of the method that was chosen.

The data will follow the flow shown in Table 1. Those steps are sequential, although not all are compulsory. In the interface, every block appears when the data has been corrected for the previous effect.

4 Background correction

Most of the published literature recommends the subtraction of the background intensity from the foreground intensity of every spot. The background intensity is defined as the intensity of the probe that attaches to the array, even when there is no cDNA available, contributing this intensity to the foreground intensity. However, the chemical properties of the array surface are still not completely known. This makes difficult to determine which is the contribution of the background intensity to the measured foreground intensity. The later will be the estimator of the abundance of mRNA for a particular gene in a particular sample.

In order to choose a suitable background correction, the interface allows the visualization of different features:

- Scatter plots. There are three different types of scatter plots that can be used for the analysis.
 - Background against foreground intensities. This plot can help to decide whether the background intensity is additive to the foreground intensity. However, since the background intensities are usually much lower than the foreground intensities, this scatter plot often suggests a linear relationship due just to the small effect that the background subtraction has in the foreground intensities.
 - Scatter plot of the background of both channels
 - Scatter plot of the foreground of both channels.

The last two will give some clues to clarify if the relationship in both, foreground and background, is similar. We could then extend conclusions from the background to the foreground and vice-versa.

- 3D plots and contour plots. These plots are useful to study the distribution of the background and foreground intensities across the array and to determine areas where the background intensity is extremely high. As shown in Figure 1 these plots gave us the chance to detect some of the gene effects. Areas around the controls shown a lower overall intensity than the rest of the array. This was due to shorter PCR products than in the rest of the array. Thanks to the possibility of plotting the contour plot of the ratio, we realized how this effect was cancelled.

All these visualization tools help to take a decision about the most suitable background correction to perform in the data set. The interface allows both possibilities: subtraction of the background or not. After choosing one of both possibilities, the scatter plots and contour plots of the corrected data set can be visualized. This can be useful to check the effect of the background correction in the data. For our data set, background subtraction appeared to pass the noise from the background to the foreground, increasing the noise in the experiment instead of reducing it. It is shown also in [Huber, von Heydebreck, Sültmann, Poustka, and Vingron \(2002\)](#) how

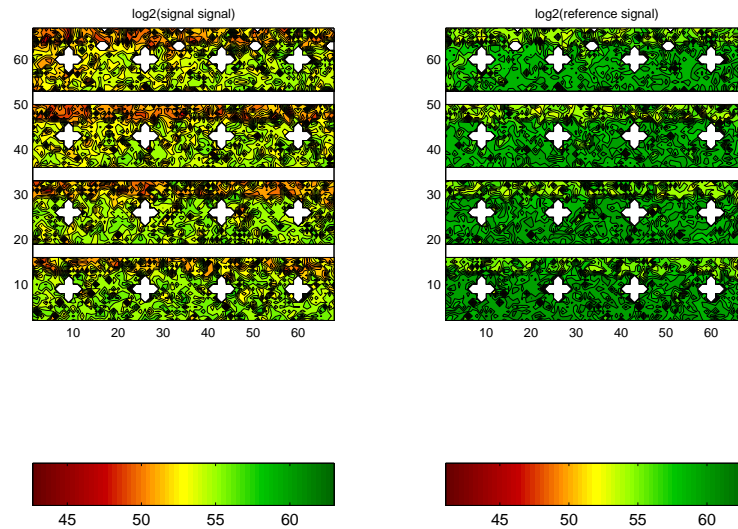


Figure 1: Contour plots for the foreground intensity of green and red channels in a particular microarray. The level curves give an idea of the distribution of the intensities according to its location in the array.

the background subtraction would increase the variability of the data using the \log_2 transformation while this variability is not so great if the background is not subtracted. Although background subtraction is currently the most popular approach, we are investigating new methods to use the information about the background intensities to correct the foreground intensities.

5 Dye correction

After background correction, systematic errors must be corrected. The most important of all of them is the one introduced by the different properties of both fluorescent dyes labelling the two RNA pools. We have detected four properties that are different for both dyes. The most important of them is the lower incorporation rate of *Cy5*, but also the quantum yield, the photobleaching and the quenching properties are different. All these differences distort the real intensity values of both channels so they must be balanced. However, we must be careful at this stage. The most popular dye correction methods are based on the idea that the majority of the genes are equally expressed in both channels. But this is not going to be the case of all experiments. For this reason, two different approaches were implemented in our interface. Using the terms defined in [Kepler et al. \(2002\)](#), the correction of the dye effect -as well known as within-array normalization- can be performed:

- using the whole data set to normalize the data (see [Figure 2](#)), as well known as self-consistency.

- using the quality control elements provided in the experiment. This includes the dye-swap normalization proposed by [Luu, Yang, Dudoit, and Speed \(2001\)](#), the use of spotted controls or the use of a reference channel (see [Figure 3](#)).

As seen in [Figure 2](#) and [Figure 3](#), both approaches can be selected in our interface.

5.1 Dye-effect correction by self-consistency

Assuming that most of the genes are going to be equally expressed in both channels, an expression ξ is estimated to force the overall intensity of both channels to be the same. Both channels intensities would be then related according to the expression:

$$R = \xi \cdot G,$$

where $R \equiv red$ and $G \equiv green$. The estimation of this expression ξ is going to result in different methods to correct the different properties of the dyes. Four of them can be selected in our interface:

Global normalization In this case, we assume that the systematical bias due to the different properties of the dyes is affecting all spotted genes in the array in the same amount. A constant k relating both channels is estimated. If most of the genes are expected to be equally expressed, then a good representative value of the distribution of the ratios is:

$$k = med_i \frac{R_i}{G_i}$$

and $\xi = k$. For experiments for which a high percentage of genes is differentially expressed comparing both channels, the use of the first or third quartiles are more suitable options. The three choices are implemented in our interface.

Linear regression normalization In [Quackenbush \(2001\)](#) a regression line is fitted to the scatter plot (G,R) . Under the assumption that most of the genes should be equally expressed for both channels, the regression line should have a slope one. Hence,

$$R = m \cdot G + n \rightarrow \frac{R}{m} - \frac{n}{m} = G .$$

From that follows $\xi \simeq m$, where m is the slope of the regression line fitted to the scatter plot and n is the intercept with the ordinate. The linear regression approach assumes that the error term has constant variance for all observations, i.e. is homocedastic. Hence, the residuals should be plotted against the independent variable G to detect possible patterns which would suggest the unsuitability of the model fitted to the data.

LOWESS normalization As suggested in [Luu et al. \(2001\)](#) and [Yang et al. \(2002\)](#), looking at the (A,M) plot implemented in our interface it can be detected if the distribution of the log ratios depends on the intensity. In this case, it is not appropriated to correct every spot in the same amount as the global method does. At the same time, the linear regression method is very sensitive to outliers, so a more robust alternative is required. For these reasons the use of a *LOWESS* function to correct the dye bias is becoming more

important in the normalization of microarray data. (A, M) scatter plot will show:

$$A_i = \frac{1}{2} \cdot (\log_2 S_i + \log_2 R_i),$$

$$M_i = \log_2 \frac{S_i}{R_i}.$$

The *LOWESS* function $c(A_i) : I \mapsto \mathbb{R}$ can be calculated from this plot, where the set of indexes I denotes all genes spotted on the array. The fitting of the *LOWESS* function $c(A)$ from the (A, M) scatterplot leads to:

$$M = \log_2 \left(\frac{R}{G} \right) \cong c(A) \Rightarrow \xi = k(A) = 2^{c(A)}.$$

To estimate this function in **MATLAB** takes extremely long, but we can improve its efficiency using a **C++** function implemented by the Jackson's laboratory (see [Churchill \(2002\)](#)). *LOWESS* is computationally efficient also in **R**.

LOWESS for different print tips During the spotting process, the spots located in the same "grid" are printed by the same print tip. [Yang et al. \(2002\)](#) suggest that different *LOWESS* functions should be fitted for the different print tip subgroups. In our interface we have implemented the scatter plots that show the genes ordered like they were spotted in the slide to detect print tip effects. The function to correct these effects if necessary is also implemented in the GUI. However, we would expect this effect to cancel with the ratios in two color-microarrays.

Regardless to the method used to estimate ξ , any of them corrects the data so,

$$\frac{R}{G} \cong 1 \Rightarrow M = \log_2 \frac{R}{G} \cong 0$$

For this reason, to look at the scatter plots, boxplots and kernel fitted functions before and after the correction is essential.

It should not be forgotten that the interface also allows the visualization of different probability plots before and after the dye correction. This feature is important in the study of the distribution of both channels intensities and this must be considered if we want to use ANOVA or any other probabilistic framework for further analysis of the data.

5.2 Dye-effect correction using the quality elements provided in the experiment

In general, there are many experiments for which the assumption of most genes equally regulated cannot be known "a priori" or for which a very different number of genes is expected to be differentially expressed in both channels. In those case we would rely on the quality control elements to normalize the data. We have implemented two methods:

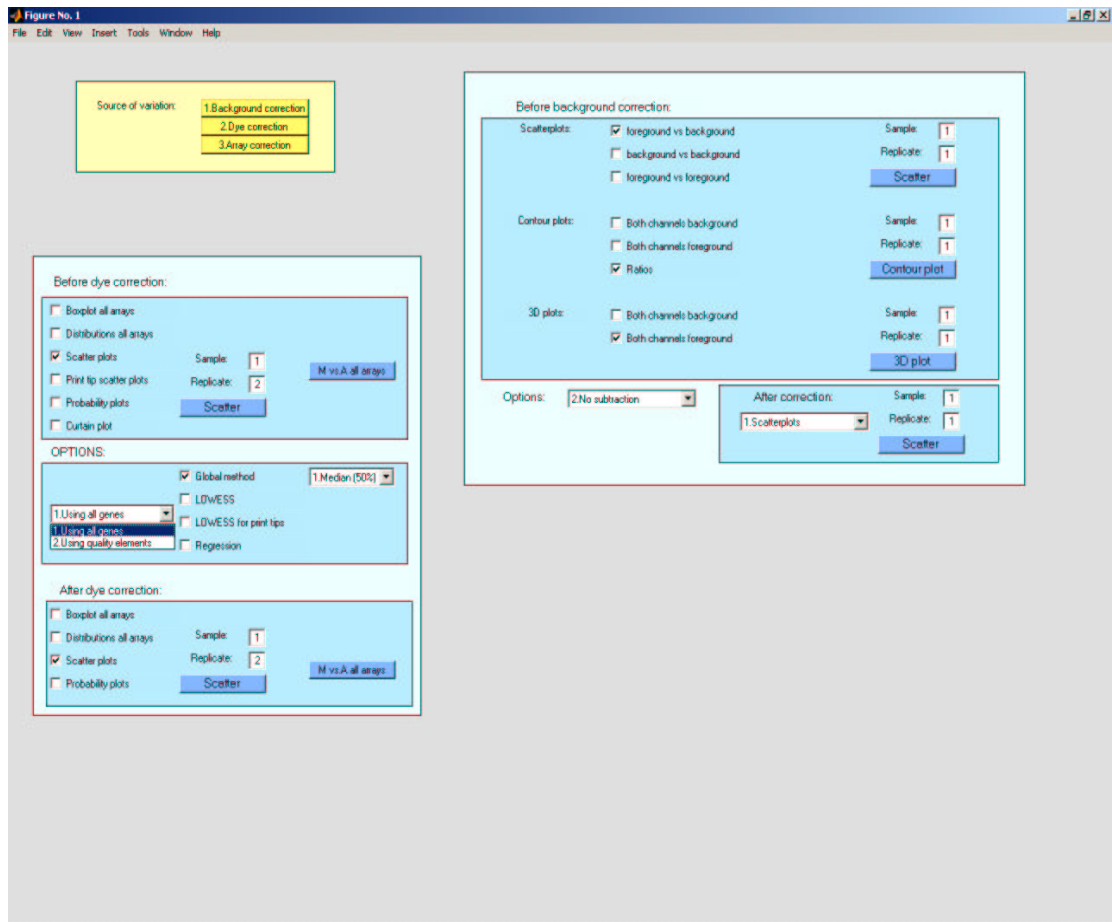


Figure 2: View of the interface after background correction and choosing the option to correct the dye effect using all the genes. We can see how all the methods outlined in Table 1 to correct the dye effect using the self-consistency approach are implemented in the interface.

Dye-swap normalization It was first described in Luu et al. (2001). Given two arrays for which the same material was labelled with a different dye each time, for every spotted gene i the following expressions are considered

$$M_i = \log_2 \left(\frac{R_i}{G_i} \right),$$

$$M'_i = \log_2 \left(\frac{R'_i}{G'_i} \right).$$

From these two equations, we obtain

$$M_i = \log_2 \left(\frac{R_i}{G_i} \right) = \log_2 \left(\frac{s_i}{r_i} \cdot k_i \right) = \log_2 \left(\frac{s_i}{r_i} \right) + \log_2 k_i = \log_2 \left(\frac{s_i}{r_i} \right) + c_i,$$

$$M'_i = \log_2 \left(\frac{R'_i}{G'_i} \right) = \log_2 \left(\frac{r_i}{s_i} \cdot k'_i \right) = -\log_2 \left(\frac{s_i}{r_i} \right) + \log_2 k'_i = -\log_2 \left(\frac{s_i}{r_i} \right) + c'_i,$$

where r_i stands for the intensity of the gene i in sample r and s_i for the same value in sample s . The target is to estimate $\log_2 \left(\frac{s_i}{r_i} \right)$ from M_i, M'_i . Hence, it follows that

$$M_i - c_i = \log_2 \left(\frac{s_i}{r_i} \right).$$

$$-M'_i + c'_i = \log_2 \left(\frac{s_i}{r_i} \right).$$

For this expression, c_i and c'_i depend on the properties of the dyes. Because they are not supposed to change significantly from one array to another it can be considered $c_i \simeq c'_i$ (see [Sanchez-Cabo, Cho, Butcher, Hinds, Trajanoski, and Wolkenhauer \(2003\)](#) for explanation). Adding both equations,

$$M_i - M'_i \simeq 2 \cdot \log_2 \left(\frac{s_i}{r_i} \right) \implies \frac{1}{2} \cdot (M_i - M'_i) \simeq \log_2 \left(\frac{s_i}{r_i} \right).$$

The main advantage of the dye-swap normalization is that transforms the data preserving the characteristics of every singular gene. Note that the computational cost for the implementation of this method is very low.

Using the controls If controls covering the whole intensity range are available, we can normalize our data using them. For controls for which the expression level in both channels is expected to be the same, a non-linear function can be fitted to the (A, M) plot of the controls and used to correct the entire data set. However, because the number of controls available per slide is usually not very large, we do not recommend to fit a *LOWESS* function but a more general method such as Levenberg-Marquardt ([Marquardt \(1963\)](#)). The model used will be in most of the cases a quadratic function.

6 Replicate handling

Besides the systematical error introduced in every measurement, there is an *error* corresponding to the random error and that cannot be perfectly estimated. The only way to reduce the intrinsic variability of a given measurement is replicating measurements. In our interface, we first consider a feature that tests the quality of the replicates. It is called curtain plot because it appears as a curtain (see [Figure 4](#)). Different correlation measurements are allowed (Standard, Pearson and Spearman). The last of them calculates the correlation in terms of the shape of gene profiles. A percentage of the genes which profiles are correlated for the different replicates is calculated. After the dye correction, the effect can be checked as well in the

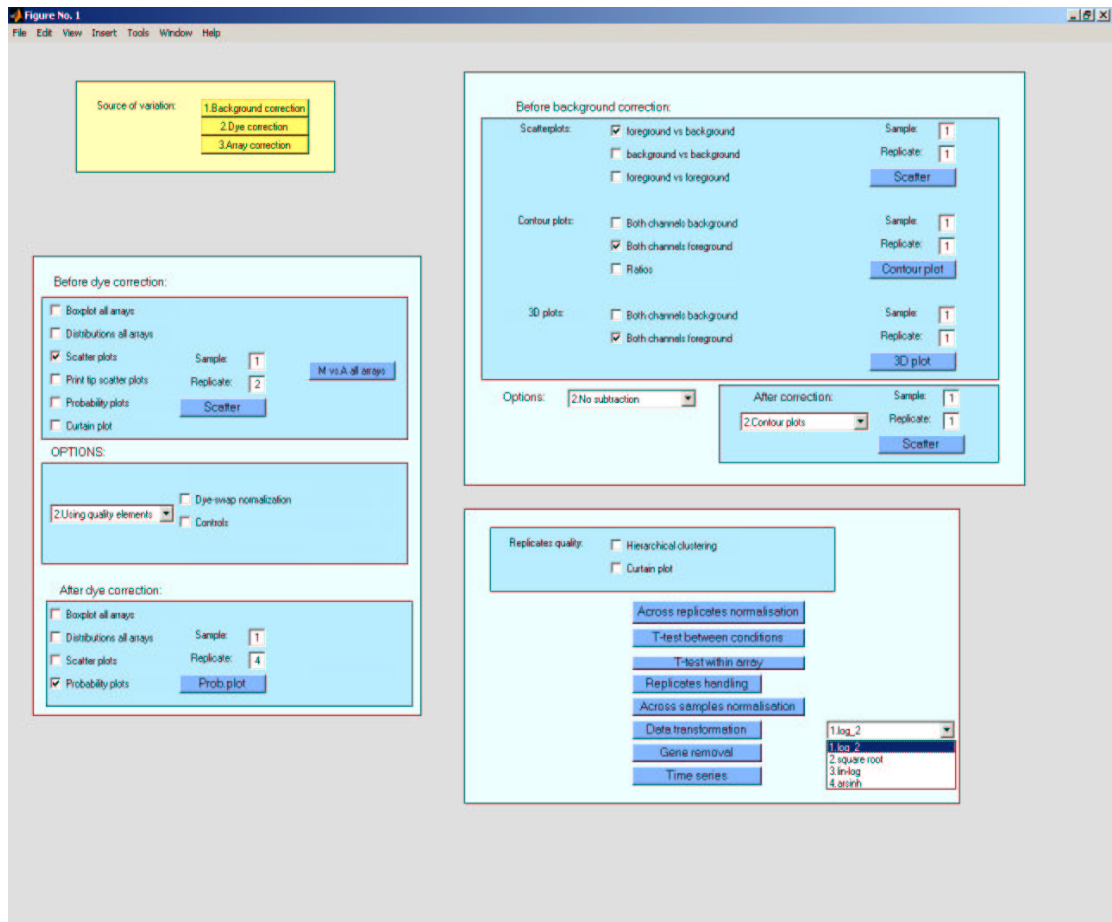


Figure 3: Interface showing the options to correct the dye-effect using the quality control elements. The different options to transform the data are as well shown.

replicates. Besides the curtain plot described before, hierarchical clustering on the replicates can be visualized.

The replicates must be used to obtain a most reliable measurement of every particular ratio. A representative value of the intensity ratio of both channels for a given spot must be taken. We choose the mean. Before taking the average, the interface allows the possibility to normalize across replicates. If the overall ratio expression level is expected to be one, the different replicates can be brought to a common reference scale dividing the different replicates for a representative value of the replicates data set, which is going to be the median:

$$t_i = \frac{q_i}{med_j(q_j)}$$

where med_j is the median of the normalized ratios q_j across the different replicates. In our interface this is called *across replicates normalisation* and contributes to reduce the experimental error that can come from inconsistent experimental conditions from array to array (see Figures 2, 3).

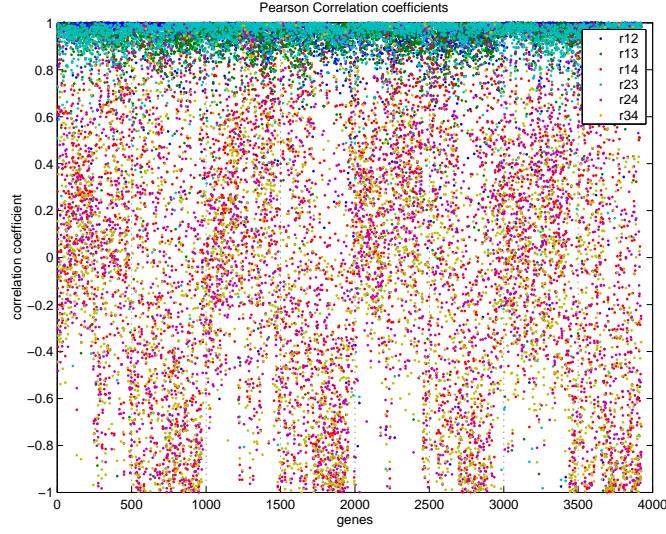


Figure 4: “Curtain” plot: The correlation coefficient for the different replicates of all the genes in the microarray are represented. It can be seen how good the quality of the replicates is. The uncorrelated replicate is the fourth, for which the dyes were swapped.

As important as getting a unique value that will be considered as the ratio of the expression level of a particular gene in both channels, is giving an estimator of how reliable this value is. For this reason we calculate as well the standard error and we use it to show in the time series plot the reliability of this value as an estimator of the log ratio.

Two different t-test are available: Equation (1) shows the formula for reference designs (Kerr and Churchill (2001)) for which the interest is in expression level between biological conditions.

$$t_{ic_1c_2} = \frac{\bar{x}_{ic_1} - \bar{x}_{ic_2}}{\sqrt{\frac{s_{ic_1}^2}{n_{c_1}} + \frac{s_{ic_2}^2}{n_{c_2}}}}, \quad (1)$$

where

$$\bar{x}_{ic_1} = \frac{1}{n_{c_1}} \sum_{j=1}^{n_{c_1}} x_{ij} = \frac{1}{n_{c_1}} \sum_{j=1}^{n_{c_1}} \log_2 \frac{R_{ij}}{G_{ij}}, \text{ and}$$

$$s_{ic_1}^2 = \frac{1}{n_{c_1} - 1} \sum_{j=1}^{n_{c_1}} (x_{ij} - \bar{x}_{ic_1})^2.$$

Equation (2) shows the formula for loop designs (Kerr and Churchill (2001)) for which the expression level of the gene is estimated using the log ratio of the intensity of the two hybridized samples.

$$t_{ic_1c_2} = \frac{\bar{x}_{ic_1c_2}}{\sqrt{\frac{s_{ic_1c_2}^2}{n_r}}}, \quad (2)$$

where

$$\bar{x}_{ic_1c_2} = \frac{1}{n_r} \sum_{j=1}^{n_r} x_{ic_1c_2j} = \frac{1}{n_r} \sum_{j=1}^{n_r} \log_2 \left(\frac{R_{ic_1}}{G_{ic_2}} \right)_j, \text{ and}$$

$$s_{ic_1c_2}^2 = \frac{1}{n_r - 1} \sum_{j=1}^{n_r} (x_{ic_1c_2j} - \bar{x}_{ic_1c_2})^2.$$

After taking the average and calculating the standard error and t-test using the available replicates per biological condition, we still need a common reference to compare the log ratio across different biological condition. For this reason *normalization across samples* is still required. We correct the data according to:

$$e_i = \frac{t_i}{med_j(t_j)}$$

where med_j is the median of the gene expression level across the different arrays. An alternative is:

$$e_i = \frac{t_i}{med_J(t_j)}$$

where med_J is the median of the gene expression level across a fixed set of arrays. It is important to remark that across replicates normalization and across samples normalization are not compulsory and are just appropriated if the overall expression from one array to another is expected to be similar. Otherwise, we would be falsely correcting the data.

The last thing previous to proper analysis of the data (e.g. clustering, time series analysis or Principal Components)¹ is the transformation that should be chosen. The \log_2 transformation is the most popular but may not be the most suitable one, due to the extreme difference between small values. For this reason a log-linear transformation or the arcsinh transformation recommended in Huber et al. (2002) can be more appropriated. All of them are implemented in our toolbox (see Figure 3).

7 Conclusions

This paper describes a GUI for the normalization of microarray data. This interface is a good example of efficient organization of one of the approaches to normalize microarray data: the sequential method. In this paper it was demonstrated the necessity of visualizing different features of the data in order to choose among different options at every stage of the sequential normalization process. Every of those steps are described in this paper. Furthermore, it is needed to test how the data changes after every different correction to avoid over-fitting or transformations that don't preserve the biological meaning of the data. The GUI for normalization allows this comparison to the user. In summary, this GUI is an example to encourage the implementation of "user friendly" environments in powerful software packages such as R. Just a few feature were included in the paper. For further information and files search <http://www.sbi.uni-rostock.de/>.

¹Implemented in MADE

References

- O. Alter, P.O. Brown, and D. Botstein. Singular value decomposition for genome-wide expression data processing and modeling. *PNAS*, 97(18):10101–10106, 2000.
- G.A. Churchill. The ANOVA project from the Jackson’s laboratory. In <http://www.jax.org/research/churchill/software/anova/Rmaanova/>, 2002.
- W.S. Cleveland. Robust locally weighted regression and smoothing scatterplots. *Journal of the American Statistical Association*, 74:829–836, 1979.
- C. Colantuoni, G. Henry, S. Zeger, and J. Pevsner. Local mean normalization of microarray element signal intensities across an array surface: Quality control and correction of spatially systematic artefacts. *Biotechniques*, 32:1316–1320, 2002.
- M.B. Eisen and P.O. Brown. DNA arrays for analysis of gene expression. *Methods Enzymol.*, 303:179–205, 1999.
- W. Huber, A. von Heydebreck, H. Sülthmann, A. Poustka, and M. Vingron. Variance stabilization applied to microarray data calibration and to the quantification of differential expression. *Bioinformatics*, 1(18):96–104, 2002.
- T.B. Kepler, L. Crosby, and K.T. Morgan. Normalization and analysis of DNA microarray data by self-consistency and local regression. *Genome Biology*, 3(7):research0037.1–0037.12, 2002.
- K. Kerr and G.A. Churchill. Experimental design for gene expression microarrays. *Biostatistics*, 2:183–201, 2001.
- K. Kerr, M. Martin, and G.A. Churchill. Analysis of variance for gene expression microarray data. *Journal of Computational Biology*, 7:819–837, 2000.
- C. Kooperberg, T.G. Fazio, Delrow J.J., and T. Tsukiyama. Improved background correction for spotted DNA microarrays. *Journal of Computational Biology*, 9(1): 55–66, 2002.
- P. Luu, Y. H. Yang, S. Dudoit, and T. P. Speed. Normalization for cDNA microarray data. In *SPIE BIOS 2001*, 2001.
- D.W. Marquardt. An algorithm for least squares-estimation of nonlinear parameters. *Journal of the Society for Industrial and Applied Mathematics*, 11:431–441, 1963.
- J. Quackenbush. Computational analysis of microarray data. *Nature Reviews Genetics*, 2(6):418–427, 2001.
- F. Sanchez-Cabo, K.Y. Cho, P. Butcher, J. Hinds, Z. Trajanoski, and O. Wolkenhauer. Is *LOWESS* a panacea in the normalization of microarray data? *Applied Bioinformatics (In press)*, 2003.
- A. Schulze and J. Downward. Navigating gene expression using microarrays - A technology review. *Nature Cell Biology*, 3:190–195, 2001.
- A. Unwin. R objects, two interfaces! (R objects to interfaces?). In Kurt Hornik and Friedrich Leisch, editors, *Proceedings of the 2nd International Workshop on Distributed Statistical Computing, March 15-17, 2001, Technische Universität Wien, Vienna, Austria*, 2001. URL <http://www.ci.tuwien.ac.at/Conferences/DSC-2001/Proceedings/>. ISSN 1609-395X.

Y.H. Yang, S. Dudoit, 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):e15.1–e15.10, 2002.

Corresponding author

Olaf Wolkenhauer
Department of Computer Science
University of Rostock
Albert Einstein Str. 21
D-18051 Rostock
Germany
Tel.: +49/381/49833-35
Fax: +49/381/49833-99
E-mail: wolkenhauer@informatik.uni-rostock.de