# Manual for arrayQCplot 2.0

## 1 Introduction

arrayQCplot is a software for the exploratory analysis of microarray data focusing on quality control. It helps biologists check the quality of their data visually and easily. This software is built on R and provides a user-friendly graphical interface for the graphics and statistical analysis. Therefore, naïve users can use arrayQCplot to check their data just by clicking mouse button.

arrayQCplot is a tool for exploring microarray data, especially two channel cDNA microarray data. It is developed to help biologists investigate microarray data in the first step of analysis. For two channel microarray data, arrayQCplot provides the normalization using lowess method, examine the reliability of experiments and check the reproducibility of replicates if available. For the other types of microarray data, for example, one channel cDNA microarray data, Affy data, etc., the limited service is provided with the normalized data.

## 2 How to install

arrayQCplot is built on the open-source statistical software, R, and provides a user-friendly interface. To use arrayQCplot package, you need to install

- R

  - a free software environment for statistical computing and graphics
  - available at http://www.r-project.org/

- GTK+

  - The premier open-source graphical user interface toolkit.
  - available at http://www.gtk.org

- RGtk2 package

  - allows one to create GTK+ and Gnome widgets from within R, using R functions and expressions
  - available at http://www.ggobi.org/rgtk2/

- cairoDevice R package

  - help create embedded graphics devices within GUIs
  - available at http://www.ggobi.org/rgtk2/

- LPE package library

  - originally developed to do significance analysis of microarray data with small number of replicates

1

- we use this package to estimate the lpe variance within treatments
- available at http://www.bioconductor.org/packages/bioc/stable/src/contrib/html/LPE.html

- arrayQCplot version 2.0

  - available at http://biostats.snu.ac.kr/ graceslee/arrayQCplot/arrayQCplot.html
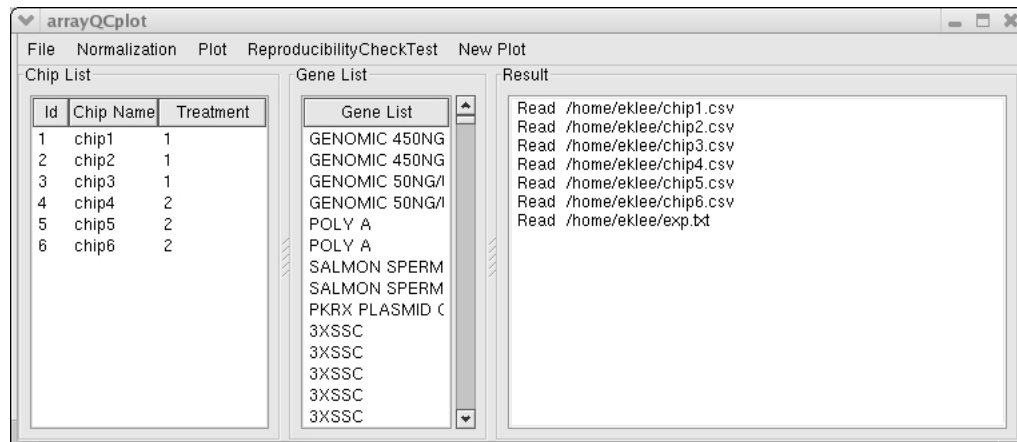
# 3    How to use arrayQCplot

- Start R

- Load the arrayQCplot library
  >library(arrayQCplot)

- Open arrayQCplot
  >arrayQC.plot()

## 3.1    Main GUI

The main GUI of arrayQCplot has the following parts :

- `Chip List` table shows the information about loaded experiments including chip name, and treatments. From this table, the user can select chips that he/she wants to analyze.
- `Gene List` shows the names of loaded genes.
- `Result` table is used to write the result of statistical analysis. In this table, arrayQCplot points out what the user did. All the result from this table can be saved as a file.
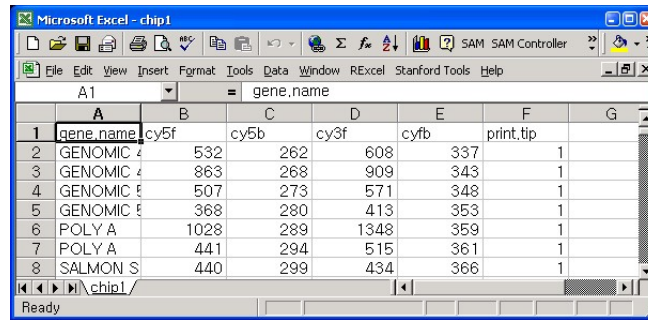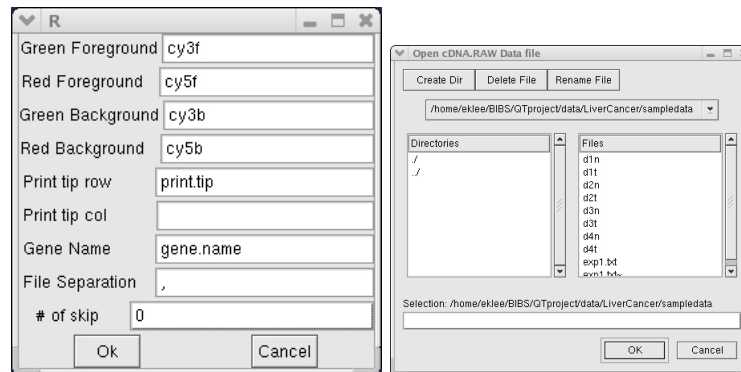
## 3.2  File

- **Read cDNA Raw data**

  To read two channel cDNA raw data, you need to save your data as following.

  - one chip should be in one file
  - In one file, arrayQCplot needs foreground intensities of cy3 and cy5, background intensities of cy3 and cy5, print tip information, and gene names.
  - recommendation : use Excel spread sheet and save your data as CSV(Comma delimited)



- **Read cDNA Normalized data**

  To read the normalized data, use this menu. In one file, the user needs to input gene names, normalized intensities for each chip.
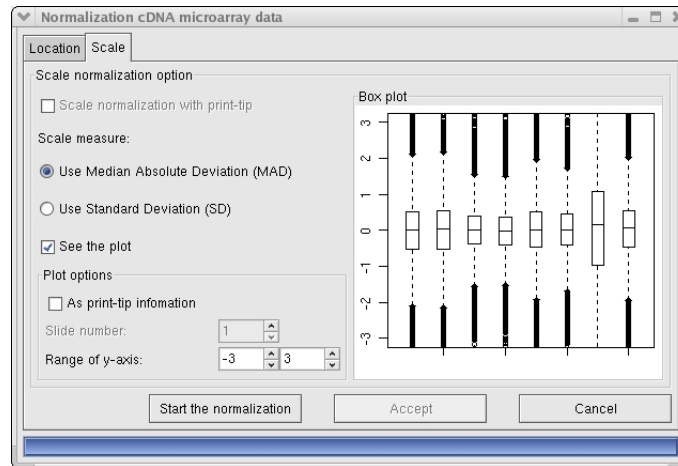
- **Read Exp.Info**

  Some of arrayQCplot functions need experimental information, especially treatment information.

- **Save Result**

  arrayQCplot can save all the results on the `Result` table.

## 3.3 Normalization

If the cDNA raw data is loaded, the user needs to normalize data before further analysis. arrayQCplot provides location and scale normalization separately. After `Location normalization`, the user can select `Scale normalization`. After all normalization procedure, the user needs to select `Accept` button. arrayQCplot provides different GUI for normalization. MA plots and box plots are also provided in the Normalization GUI.
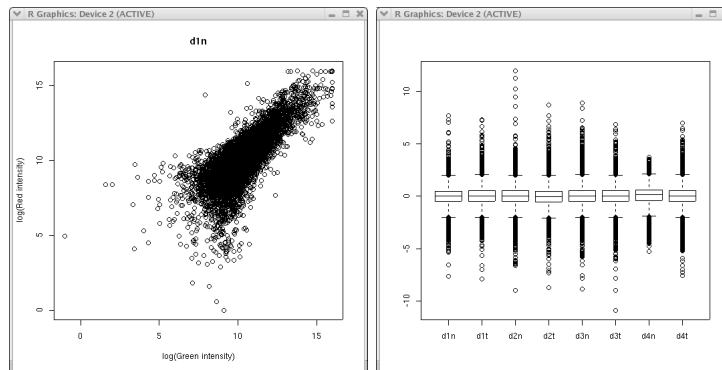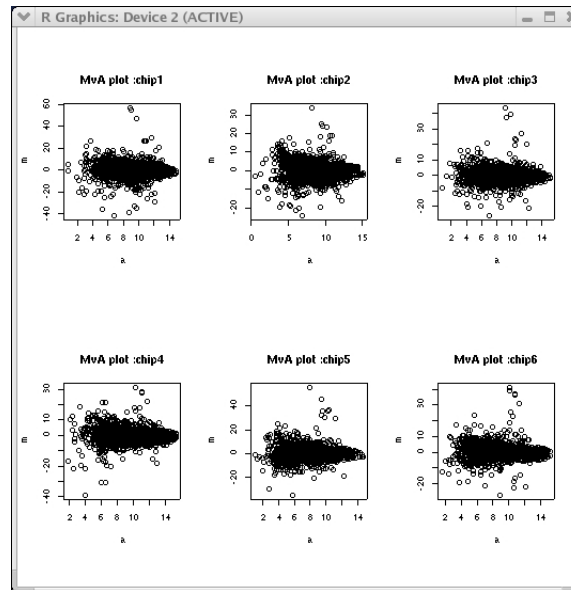
## 3.4 Plot

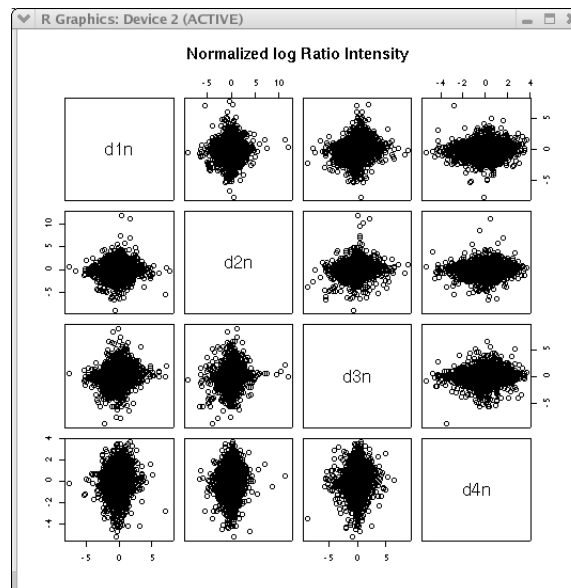For the normalized data, arrayQCplot provides a variety of basic plots to explore microarray data.
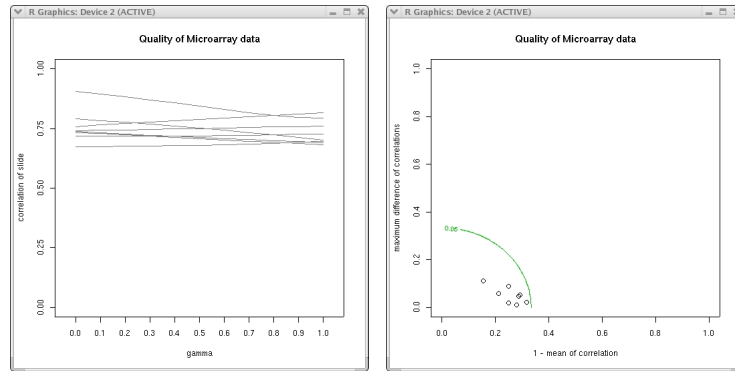
- **Red vs. Green plot**
- **Box plot**
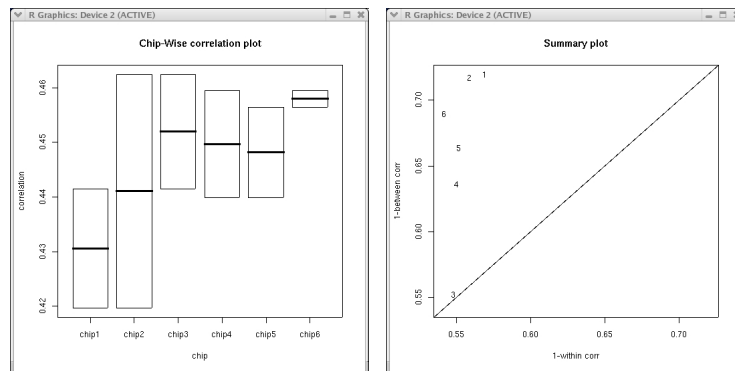
- **MvA plot**



- **Scatter plot matrix : Red/Green/Ratio**

- **Quality Control plot : correlation plot/diagnostic plot**
  To detect outlying chips, arrayQCplot provides two plots, correlation plot and diagnostic plot. For more detailed explanation, see Park et al.(2005).



- **Chip-wise correlation plot : Summary correlation plot**
  To check the reproducibility and quality of experiments, arrayQCplot provides two plots, chip-wise correlation plot and summary correlation plot. For more detailed explanation, see Lee et al.(2006).

## 3.5 Reproducibility Check Test

**Reproducibility Check Test** menu in arrayQCplot focuses on the reproducibility, sensitivity and specificity. Therefore, to use this menu, the experiment should have at least two replicates for each treatment level. Let

$Y_{ijg}$ : the normalized log intensity of the $j$th replicates of the $i$th treatment for gene $g$

$\quad i = 1, 2, \cdots, I, \quad j = 1, 2, \cdots, n_i, \quad \text{and} \quad g = 1, 2, \cdots, G$
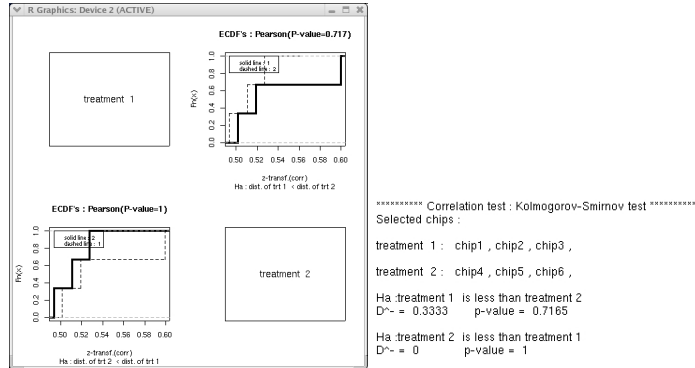
$Y_{ij} = [Y_{ij1}, Y_{ij2}, \cdots, Y_{ijG}]^T$

$r_{ij,kl} = corr(Y_{ij}, Y_{kl})$

$\mathbf{R}_i^w = \{r_{ij,il}, \forall l \neq j\}$
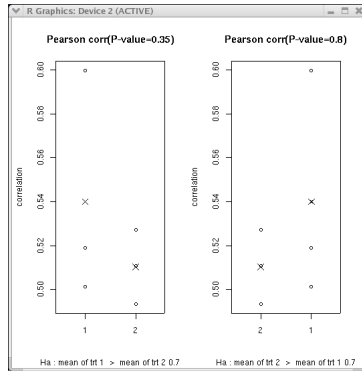
- **Test based on correlation**

  1. **K-S test with z-transformation**
     To compare the distributions of correlation sets of two treatments, $\mathbf{R}_i^w$ and $\mathbf{R}_{i'}^w$, we use Kolmogorov-Smirnov test after z-transformation. If $p - value$ is small, $\mathbf{R}_i^w$ and $\mathbf{R}_{i'}^w$ have different distributions. If the empirical CDF's of $\mathbf{R}_i^w$ is lower than the empirical CDF's of $\mathbf{R}_{i'}^w$, the reproducibility and sensitivity of treatment $i$ is better than treatment $i'$.



  2. **mean difference test**
     To test the differences of correlation sets of two treatments, $\mathbf{R}_i^w$ and $\mathbf{R}_{i'}^w$, we use Wilcox rank sum test. If $p - value$ is small, $\mathbf{R}_i^w$ and $\mathbf{R}_{i'}^w$ are significantly different. If the mean of $\mathbf{R}_i^w$ is greater than the mean of $\mathbf{R}_{i'}^w$, the reproducibility and sensitivity of treatment $i$ is better than treatment $i'$.
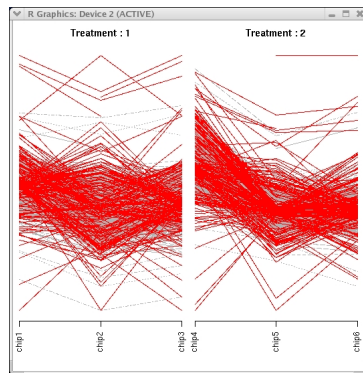
8

- ## Test based on intensities

  ### Within treatment

  For gene $j$, we test whether the mean intensities within treatment are same or not. We assume that for gene $g$, $Y_{ijg} \sim N(\mu_{ijg}, \sigma_{ig}^2)$.

  Test for each treatment : To check the reproducibility of the treatment $i$, we test $H_0 : \mu_{i1g} = \mu_{i2g} = \cdots = \mu_{in_ig}$. After applying step 1 in the overall test, $\quad \sum_j \left( \frac{Y_{ijg} - \bar{Y}_{i.g}}{\hat{\sigma}_i} \right)^2 \sim \chi^2(n_i - 1)$.



# 4    References

Park, T., Yi, S-G., Lee, S. Y., and Lee, J. K. (2005) Diagnostic plots for detecting outlying slides in a cDNA microarray experiment. *Biotechniques 38:463-471*
Lee, E. K., Yi, S-G, and Park, T. (2006) Exploratory Methods for Checking Quality of Microarray data. *Technical report, BIBS, Seoul National University*