

# R/MAANOVA: An extensive R environment for the Analysis of Microarray Experiments

Hao Wu, Hyuna Yang, Gary A. Churchill

March 8, 2008

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Installation</b>	<b>3</b>
2.1	System requirements . . . . .	3
2.2	Obtaining R . . . . .	3
2.3	Installation - Windows(9x/NT/2000) . . . . .	3
2.4	Installation - Linux/Unix . . . . .	3
<b>3</b>	<b>Data and function list</b>	<b>5</b>
<b>4</b>	<b>Preparing the input files</b>	<b>7</b>
4.1	Preparing the data files . . . . .	7
4.2	Preparing the design file . . . . .	8
4.3	Preparing the covariate matrix . . . . .	9
<b>5</b>	<b>A quick tour of the functions</b>	<b>10</b>
5.1	CAMDA kidney experiment . . . . .	10
5.2	An Affymetrix experiment . . . . .	20
<b>A</b>	<b>Running R/maanova on computer cluster</b>	<b>25</b>
A.1	System requirements . . . . .	25
A.2	System setup . . . . .	25
A.3	Install and test clusters in R . . . . .	26
<b>B</b>	<b>Basic Algorithms</b>	<b>27</b>
B.1	ANOVA model for microarray experiment . . . . .	27
B.1.1	Model . . . . .	27
B.1.2	Fixed model versus Mixed model . . . . .	27

B.2 Hypothesis tests . . . . .	28
B.2.1 F test in matrix notation . . . . .	28
B.2.2 Building $L$ matrix . . . . .	29
B.2.3 Four flavors of F-tests . . . . .	29
B.3 Data shuffling in permutation test . . . . .	30
<b>C Frequently Asked Questions</b>	<b>31</b>

# 1 Introduction

*R/maanova* is an extensible, interactive environment for the analysis of microarray experiments. It is implemented as an add-on package for the freely available statistical language R ([www.r-project.org](http://www.r-project.org)). MAANOVA stands for MicroArray ANalysis Of VAriance. It provides a complete work flow for microarray data analysis including:

- Data quality checks - visualization and transformation
- ANOVA model fitting - both fixed and mixed effects models
- Calculating test statistics - F and Fs statistics
- Deriving p-value - using sampling and residual shuffling permutation approach
- Summarize the result - provide summary tables, graphics including volcano plot and cluster analysis with bootstrapping

*R/maanova* can be applied to any microarray data but it is specially tailored for multiple factor experimental designs. Mixed effects models are implemented to estimate variance components and perform F and T tests for differential expressions.

## 2 Installation

### 2.1 System requirements

This package was developed under *R 1.8.0* and updated under *R 2.5.1* in the *Linux Redhat 8.0* operating system. The programs have been observed to work under *Windows NT/98/2000*. The memory requirement depends on the size of the input data but a minimum of 256Mb memory is recommended.

### 2.2 Obtaining R

R is available in the Comprehensive R Archive Network (CRAN). Visit <http://cran.r-project.org> or a local mirror site. Source code is available for UNIX/LINUX, and binaries are available for Windows, MacOS and many versions of Linux.

### 2.3 Installation - Windows(9x/NT/2000)

- Install *R/maanova* from Rgui
  1. Start Rgui
  2. Select Menu Packages, click Install package from local zip file. Choose the file `maanova_*.tar.gz` and click 'OK'.
- Install *R/maanova* outside Rgui
  1. Unzip the `maanova_*.tar.gz` file into the directory `$RHOME/library` (`$RHOME` is something like `c:/Program Files/R/rw1081`). Note that this should create a directory `$RHOME/library/maanova` containing the R source code and the compiled dlls.
  2. Start Rgui.
  3. Type `link.html.help()` to get the help files for the *maanova* package added to the help indices.

### 2.4 Installation - Linux/Unix

1. Go into the directory containing `maanova_*.tar.gz`.
2. Type R CMD INSTALL *maanova* to have the package installed in the standard location such like `/usr/lib/R/library`. You will have to be the superuser to do this. As a normal user, you can install the package in your own local directory. To do this, type R CMD INSTALL `-library=$LOCALRLIB maanova_*.tar.gz`, where

`$LOCALRLIB` is something like `/home/user/Rlib/`. Then you will need to create a file `.Renvirom` in your home directory to contain the line `R_LIBS=/home/user/Rlib` so that R will know to search for packages in that directory.

### 3 Data and function list

The following is a list of the available functions. For more information about the function usage use the online help by typing “?functionname” in *R* environment or typing `help.start()` to get the html help in a web browser.

1. Sample data available with the package. Both include subset of original data, which can be found at <http://www.jax.org/staff/churchill/labsite/software/>.

`abf1` A 18-array Affymetrix experiment.

`kidney` A 6-array kidney data set from CAMDA (Critical Assessment of Microarray Data Analysis).

2. File I/O

`read.madata()` Read microarray data from TAB delimited text file (typical input for two-color array) or matrix type *R* object (such as output of `exprs()` from normalization package, such as *Affy*, *beadarray*).

`write.madata()` Write microarray data to a TAB delimited text file.

3. Data quality check - for two-color arrays

`arrayview()` View the layout of the arrays

`riplot()` Ratio intensity plot for arrays

`gridcheck()` Plot grid-by-grid data comparison for arrays

`dyeswapfilter()` Flag the bad spots in dye swap experiment

4. Data transformation

`transform.madata()` Data transformation with options to use any of several methods

5. ANOVA model fitting

`fitmaanova()` Fit ANOVA model

`resiplot()` Residual plot on a given ANOVA model

6. Hypothesis testing

`matest()` F-test or T-test with permutation

`adjPval()` Calculate FDR adjusted P values given the result of `matest()`

7. Summarize the result

`summarytable()` Provide summary of `matest()`, such as P-value, adjusted P-value and fold change

`volcano()` Volcano plot for summarizing F or T test results

## 8. Clustering

`macluster()` Bootstrap clustering.

`consensus()` Build consensus tree out of bootstrap cluster result

`fom()` Use figure of merit to determine the number of groups in K-means cluster

`geneprofile()` Plot the estimated relative expression for a given list of genes

## 9. Utility functions

`fill.missing()` Fill in missing data.

`PairContrast()` Make all possible pairwise contrast.

`summary.madata()` Summarize the data object.

`subset.madata()` Subsetting the data objects.

## 4 Preparing the input files

A data file and a design file should be prepared manually.

### 4.1 Preparing the data files

There should be only ONE data file. Most gridding software produces one file for each slide. Thus especially when you use two-color array and have multiple files, you have to combine these files to get one data file as a input for *R/maanova*.

**Input data file** Data file can be two kinds of type.

- matrix type *R* object : *R/maanova* can read matrix *R* object including the output of `exprs()` from *affy* or *beadarray* package. `Rowname` of the matrix is used as probe (clone) ID, and it is assumed that the data is started from the first column.
- TAB delimited text file : This can be a typical input file for two-color array. *R/maanova* assumes that first column is probe (clone) ID, and intensity data is saved from the second column (i.e. `probeid = 1`, `intensity = 2`). Otherwise user needs to specify the column number storing probe ID and column number from which data starts by using `probeid` and `intensity`, respectively. When clone ID is not provided, 1,2,... are used as a probe ID.

**Input data file for N-day array** For N-dye arrays, the N channel intensity data for one array needs to be adjacent to each other (in consecutive columns). The order (dye1, dye2, ...) must be consistent across all of the slides. You can put the spot flag as a column after intensity data for each array. (Note that if you have flag, you will have N+1 columns of data for each array. Again, this must be consistent for all arrays in the data set.) If you have duplicated spots within one array, replicated measurements of the same clone on the same array should appear in adjacent rows. This can be easily done by sorting on `probeid`. The number of replicates must be constant for all genes. As an example, you have four slides for a two-dye array experiment scanned by GenePix. Then you will have four output files. Following the steps to create your data file:

1. Open your favorite spread sheet editor, e.g., MS Excel and create a new file.
2. Paste your clone ID, Gene names, Cluster ID and whatever information you want to keep into the first several columns.
3. Open your first GenePix file in another window, copy the grid location into next 4 columns (you only need to do this once because they are all the same for four slides).
4. Copy the two columns of foreground mean value (if you want to use it) and one column of flag to the file in the order of Cy5, Cy3, flag.

5. Open your other 3 files and repeat step 4.
6. Select the whole file and sort the row according to Clone ID.
7. Save the file as tab delimited text file and you are done.

CAMDA kidney data file at <http://www.jax.org/staff/churchill/labsite/software> can provide a good example.

**No missing in input** The data file must be “full”, that is, all rows have to have the same number of columns. Also data should not have any missing value. If so, you want to either check the experiment to get the data or use `fill.missing()` function to fill in the missing value.

Sometimes reading and trailing TAB in the text file can cause problems, depending on the operating system. So the you should be careful about that. Sometimes the special characters in gene description can cause reading problem. We don't encourage you to put the gene description in the data file. If you have to do that, you must be careful (sometimes you need to remove the special characters manually).

## 4.2 Preparing the design file

**Design file** Design file can be three kinds of type.

- TAB delimited text file : should be saved at the working directory.
- `data.frame` *R* object
- `matrix` *R* object.

**Size of design file** The number of row in design file should be the number of array times *N* (the number of dyes) (plus one for column header, if the design file is TAB delimited text file, and header = T). The number of columns in this file depends on the experimental design. You can specify all information related to experiment at the `designfile`, but do not necessarily need to use all of them to specify your model. Among factors in `designfile`, you can specify only some of them in a formula at `fitmaanova()`.

**Required fields in design file** You must have an **Array** column in the design file (column headers are case sensitive). For two-color (or multi-color) array, user needs two more extra columns, **Dye** and **Sample**. **Dye** is reserved for the dye name, and **Sample** is reserved to identify biological replicates and reference samples. Usually you should assign each biological individual a unique sample number. Reference samples are represented by zero(0). Reference sample are treated differently. They will always be treated as fixed factor in the model and not involved in statistical tests.

**Sample** is not an required field in one-color array, however if the experiment has technical replicate, **Sample** field should be included to get a valid result. Refer to Section 5.2 for an example.



**Sample** column in design file need to be continuous integer. All other columns can be either integers or characters.

**Names must not used in design file** You must not have following columns in the design file:

- **Spot**: reserved for spot effect
- **Label**: reserved for labeling effect
- **covM**: reserved for covariate matrix

### 4.3 Preparing the covariate matrix

**Array specific covariate** If you want to include array specific covariate, it should be included in the design file and specified in the `fitmaanova()`. The covariate variable should be an integer and the size should equal to the number of arrays.

**Gene specific covariate** If you have a gene specific covariate variable, you need to prepare the covariate matrix, `covM`, and include it at the `read.madata()` and the formula in `fitmaanova()`. `covM` can be either `matrix` R object or TAB delimited tex file. When it is a delimited file, `covM` should have column name (if `header = T` ) but no row name, and the covariate variable size should equal to the size of `datafile`. Currently this only works for the fixed model.

Refer to `help(read.madata)` and `help(fitmaanova)` for more information.

## 5 A quick tour of the functions

This section will go through a few demo scripts distributed with the package to help users understand the function syntax and capabilities. The original data file can be downloaded from:

<http://www.jax.org/staff/churchill/labsite/software/>

### 5.1 CAMDA kidney experiment

This is the kidney data from CAMDA (Critical Assessment of Microarray Data Analysis) originally described by Pritchard *et al.* The web-site for CAMDA is <http://www.camda.duke.edu>. It is a 24-array double reference design. Six samples are compared to a reference with dye swapped and all arrays are duplicated. Flag for bad spots is included in the data.

Due to the limited size of vignette, data contains only first 300 genes, and you may not reproduce the figures presented in the following sections. You can find the original data file (in text format) at <http://www.jax.org/staff/churchill/labsite/software/>.

**Read Data** Suppose you download the data and design file, then this is how you read the data.

```
R> kidney <- read.madata(datafile="kidney.txt",designfile="kidneydesign.txt",
  arrayType='twoColor', metarow=1, metacol=2, col=3, row=4, probeid=5, intensity=7)
```

**Data checking and transformation** For two-color array, you need to check if the data needs transformation. To do that you can try following data checking.

1. First load data into the workspace

```
R> data(kidney)
```

2. Then we do some data quality check

```
R> gridcheck(kidney.raw)
R> graphics.off()
```

`gridcheck` is used to check the hybridization and gridding quality within and cross arrays. User should see near linear scatter plots in all grid for all arrays. This one is not great but acceptable. The grid check plot for the first array is shown in Figure 1. The red dots are for the spots with flags.

```
R> riplot(kidney.raw)
```

`riplot` stands for ratio-intensity plot. It is also called MA plot. The riplot for the first array is shown in Figure 2.

```
R> arrayview(kidney.raw)
```

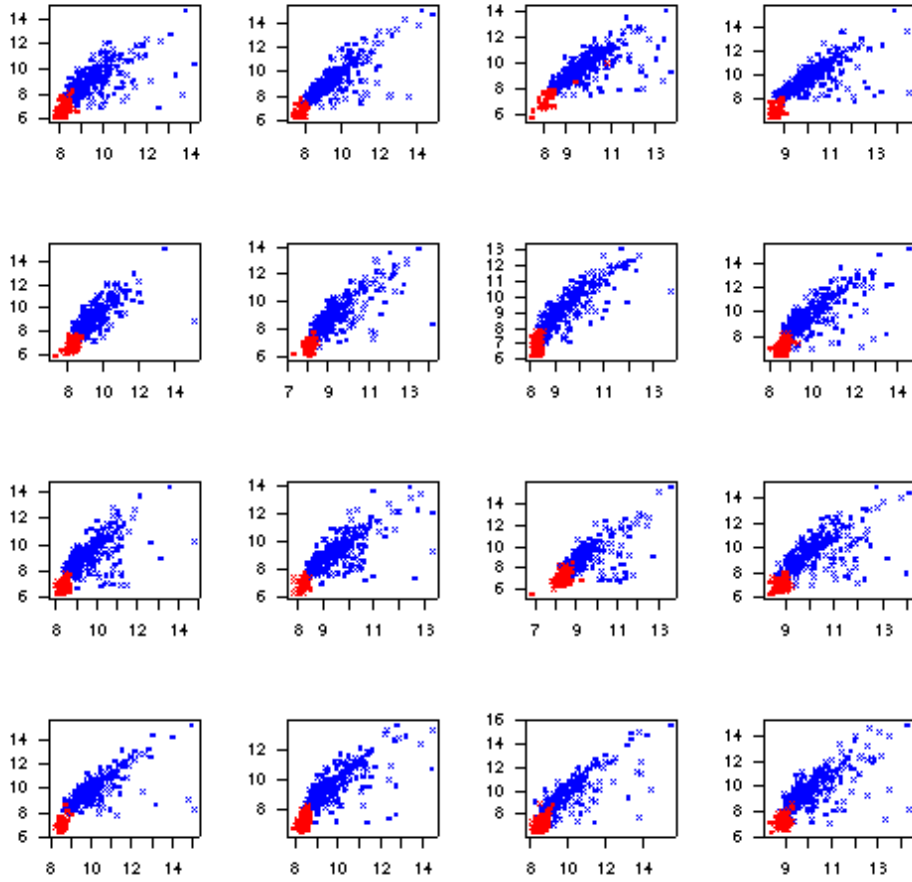


Figure 1: Grid check plot for the first array in kidney data

RI plot for array number 1

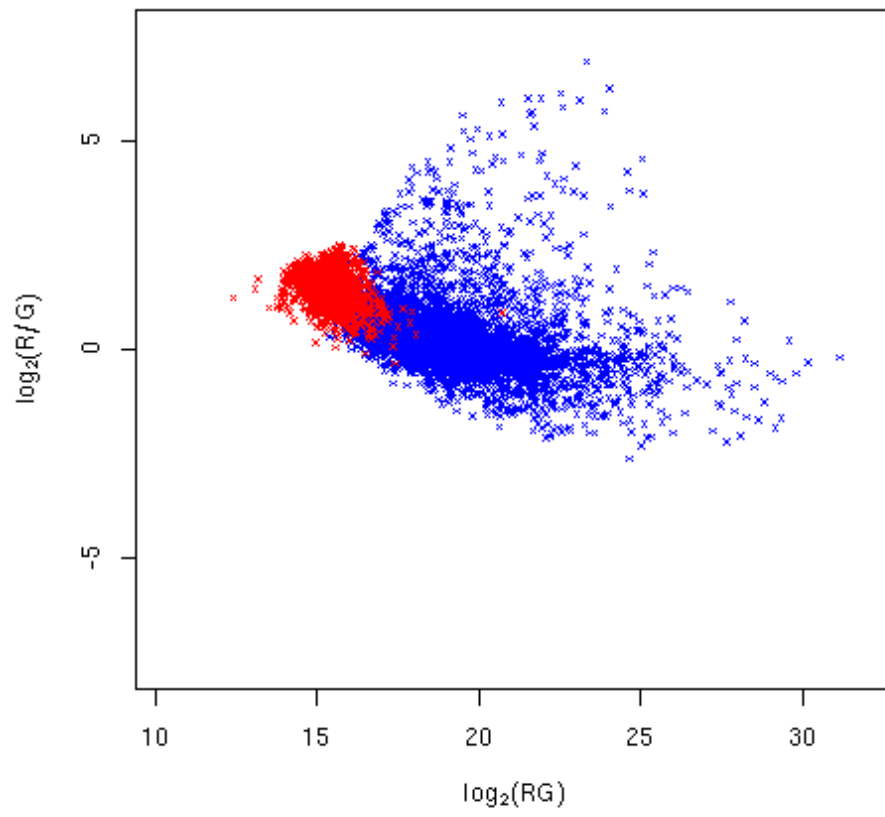


Figure 2: RI plot for the first array in kidney data

`arrayview` is used to view the spatial pattern of the arrays. The standardized log ratios for all spots are shown in different colors. The arrayview for the first array is shown in Figure 3.

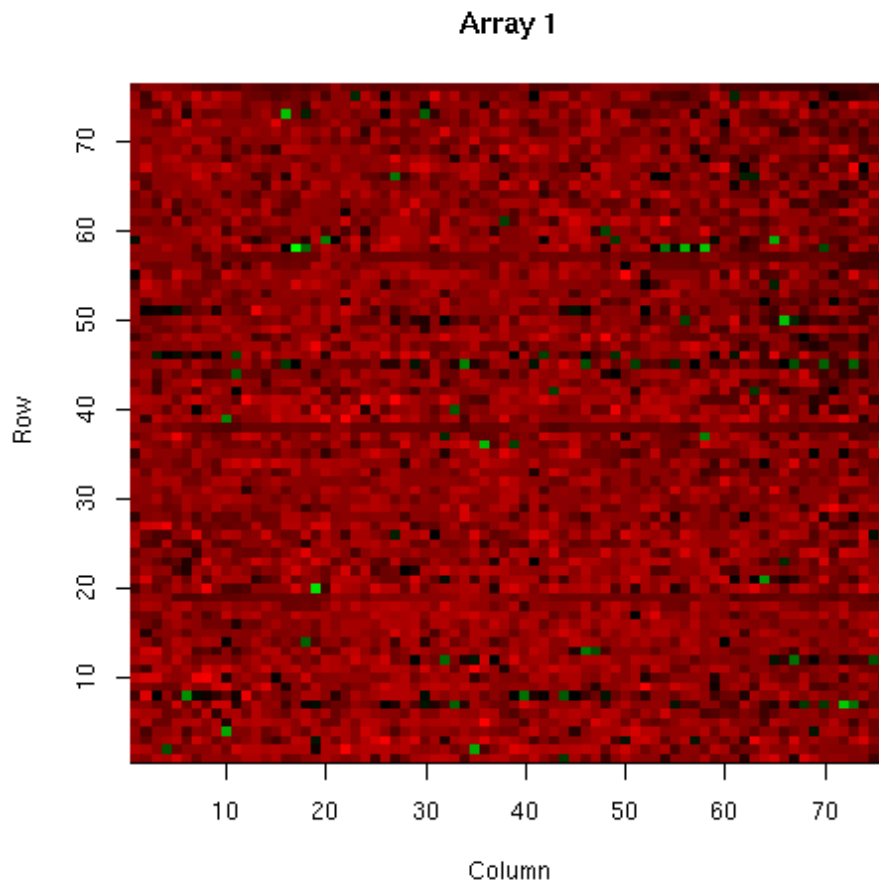


Figure 3: Arrayview for the first array in kidney data

You will generate a lot of figures by doing `gridcheck`, `riplot` and `arrayview`. Use `graphics.off()` to close all figures.

3. Transform the data using spatial-intensity joint loess.

```
R> kidney <- transform.madata(kidney.raw, method="rloess")
```

There are several data transformation method. Which method to use depends on the data. Read Cui *et al.*(2002) for detail. The transformation plot for the first array is shown in Figure 4.

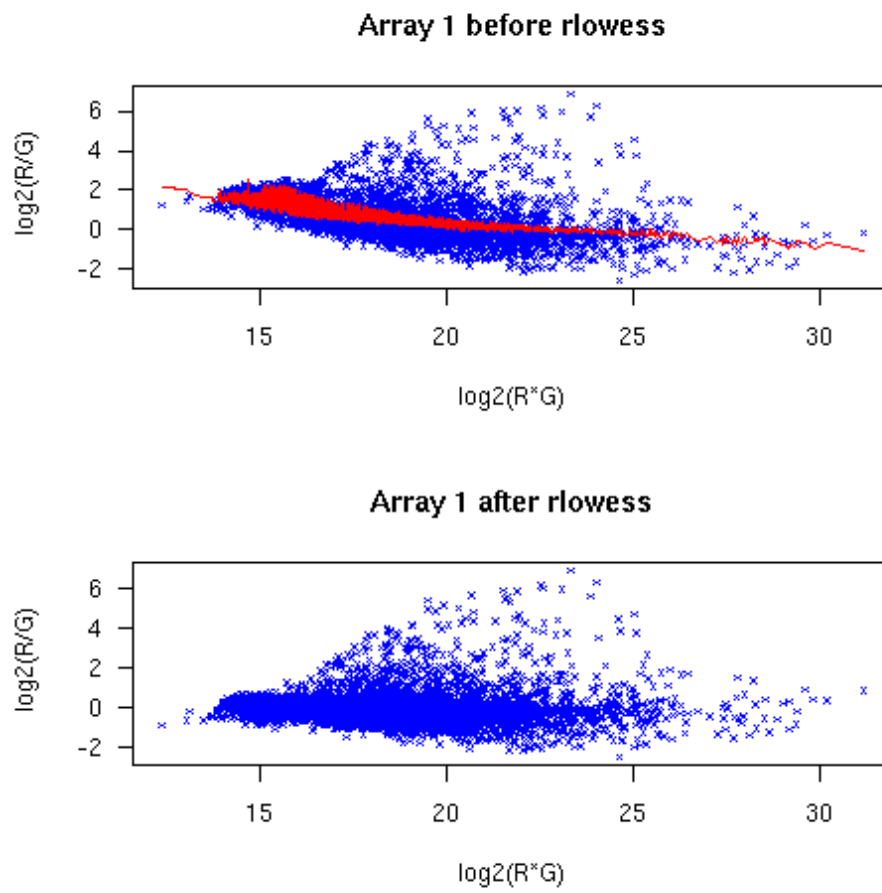


Figure 4: Joint lowess transformation on the first array for the kidney data

**Model fitting** To fit model, you need to specify `formula`, as well as `random` and `covariate` terms if one has. `formula` includes factors affecting the gene expression among ones specified in `designfile`. This also can include interactions between upto two terms. `random` and `covariate` terms should be included in `formula`. `resiplot()` shows the residual plot after model being fitted.

```
R> fit.fix <- fitmaanova(kidney, formula = ~Dye+Array+Sample)
R> resiplot(kidney, fit.fix)
```

**Test statistics** To test how significantly a given term or terms affect to the gene expression level, *R/maanova* calculates F, Fs and Fss test statistics (Fss is for fixed model only at this stage). You can include test statistics using `test.method`, in which the order is F, Fs and Fss test statistics and 1 indicates including the corresponding test statistics. As a default, *R/maanova* calculates F and Fs test statistics, and we recommend to use Fs test statistics.

*R/maanova* provides permutation method to calculate the significance of each test statistics. It implements residual shuffling and sample shuffling, and sample shuffling is the current default. Permutation test could be very time consuming, especially for the mixed effects models. The permutation test function can run on the Linux clusters through the message passing interface (MPI). For the detail information regarding F test and using computer cluster, please read appendix.

Test Sample effect in the model:

```
R> test.fix <- matest(kidney, fit.fix, term="Sample", n.perm=500)
```

FDR adjust p-value can be obtained using `adjPval`. It has four options; `stepup` (Hochberg and Benjamini, 1990), `adaptive` (Benjamini and Hochberg, 2000), `stepdown` (Westfall and Young, 1993) and `jsFDR` (Storey, 2002). `jsFDR` option uses *qvalue* package by John Storey and the package should be installed before this option is used. FDR adjusted P values:

```
R> test.fix <- adjPval(test.fix, method = 'jsFDR')
```

**Summarize the result** Each p-value or multiple testing adjusted p-value can be found under `matest()` object, i.e., p-value for Fs test statistics is saved at `test.fix$Fs$Pvalperm` and FDR adjusted p-value for Fs test statistics is saved at `test.fix$Fs$adjPvalperm`. `summarytable()` assembles fold-change, p-value and adjusted p-value (if it is available) and the result is saved in 'summarytable.csv' (default file name) under the working directory. User also can save subset of result using certain threshold.

```
R> result = summarytable(test.fix)
R> summarytable(test.fix, method=c("Pvalperm"), test=c("F1", "Fs"),
whichTest=c("Fs.Pvalperm"), threshold=.1, outfile='shrotssummarytable.csv')
```

After the F-test result is obtained, We can do volcano plot to visualize it. `volcano()` function has options to choose the P-values to use and set up thresholds. We will use tabulated P values for F and FDR adjusted permutation P values for other tests here. In the plot, the red dots above the horizontal line represent the significant genes. Note that the the flagged spot can be highlighted in the plot by providing `highlight.flag=TRUE`.

```
R> idx.fix <- volcano(test.fix,method=c("unadj", "fdrperm"),  
  highlight.flag=FALSE)
```

The volcano plot is shown in Figure 5. Note that the return variable of `volcano`

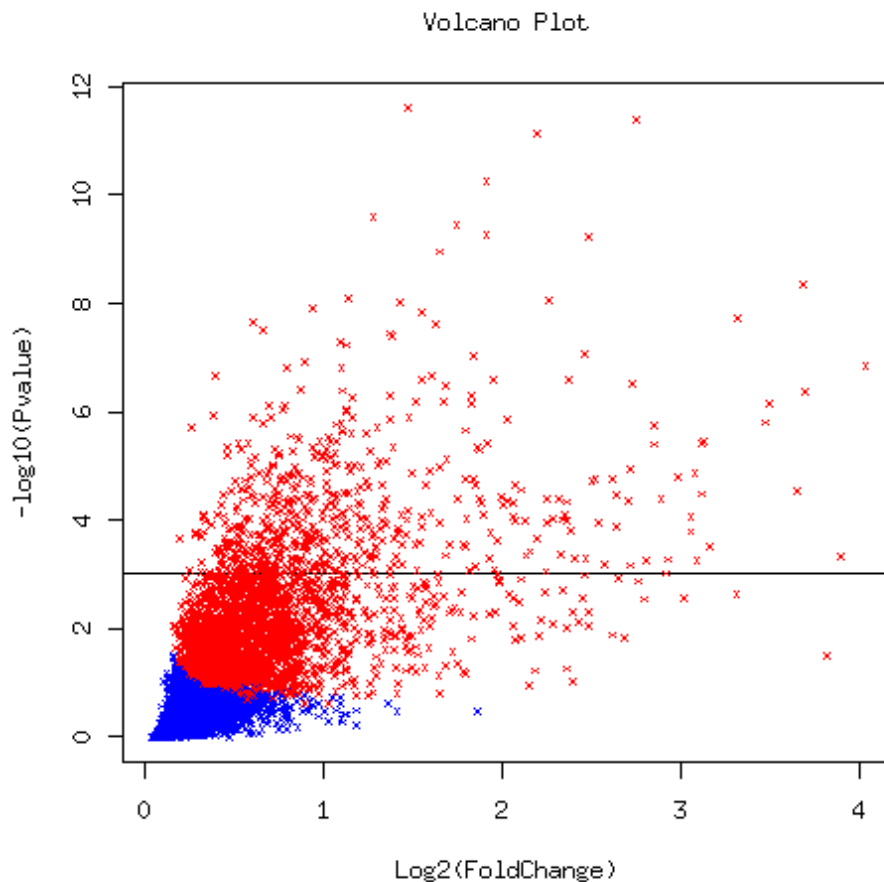


Figure 5: Volcano plot for kidney data - fixed effect model

contains the indices for significant genes.

Now we can do cluster bootstrapping and build consensus trees. Currently two cluster methods are implemented; K-mean and hierarchical clustering. Hierarchical cluster could be very sensitive to bootstrap if you have too many leaves on the cluster. Some small disturbance on the data could change the whole tree structure. Suppose you have many



genes, say, more than 50, and you want to build a consensus tree from 100 bootstrapped hierarchical trees. Then it is very likely that you will get a comb, that is, all leaves are directly under root. But if you only have a few genes to cluster, it is working fine. So we recommend to use K-means to cluster the genes and use hierarchical to cluster the samples.

`macluster` is the function to do cluster bootstrapping and `consensus` is used to build consensus trees (groups for K-means) from the bootstrap results.

```
R> cluster.kmean <- macluster(fit.fix, ,term="Sample",
  idx.gene=idx.fix$idx.all,what="gene", method="kmean",
  kmean.ngroups=5, n.perm=100)
R> con.kmean <- consensus(cluster.kmean, 0.7)
R> con.kmean$groupname
```

An expression profile plot will be generated for the consensus K-means result. The plot is shown in Figure 6, and `groupname` shows the name of genes belonging to each group.

Now we can do hierarchical cluster on the samples. The consensus tree is shown in Figure 7.

```
R> cluster.hc <- macluster(fit.fix, term="Sample",
  idx.gene=idx.fix$idx.all,what="sample",
  method="hc", n.perm=100)
R> con.hc <- consensus(cluster.hc)
```

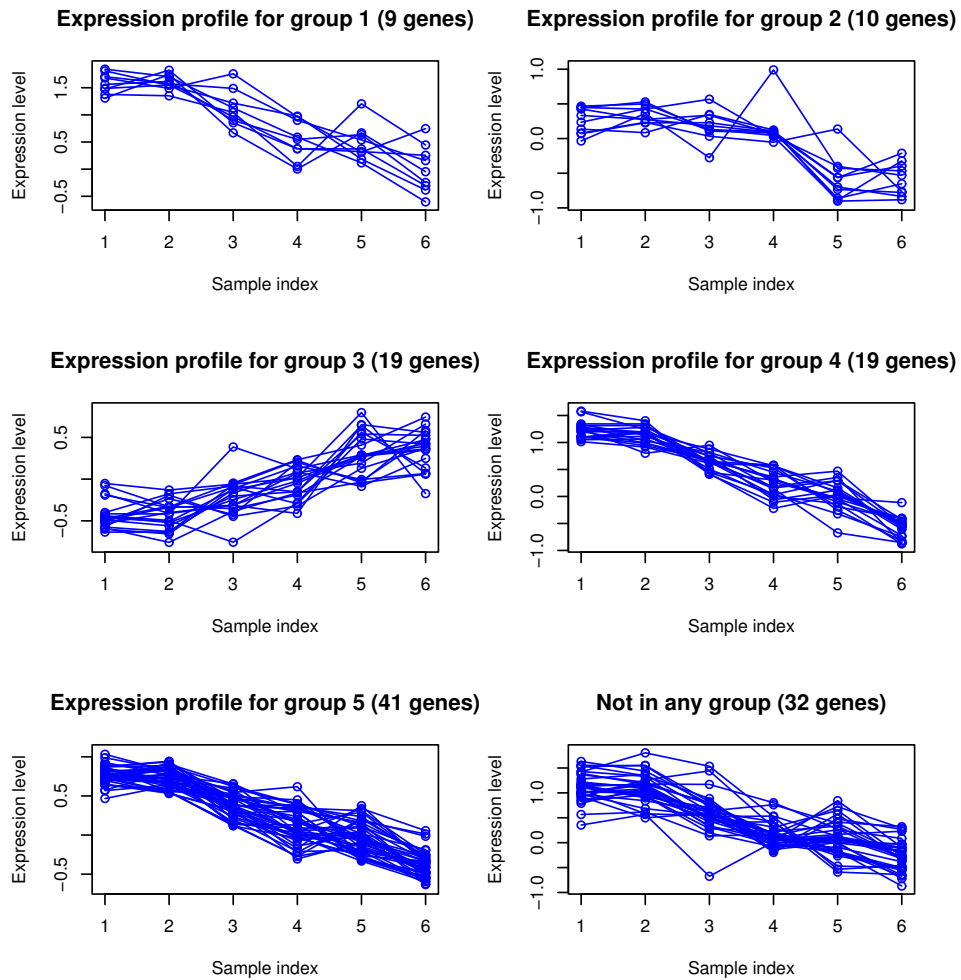


Figure 6: Expression profile plot for bootstrap K-means result, kidney data

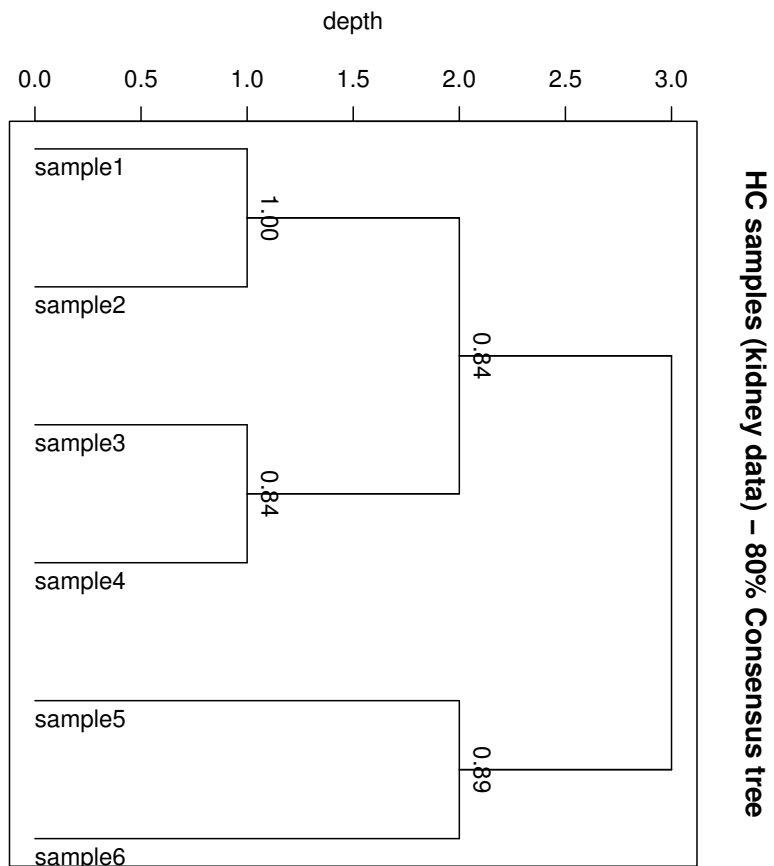


Figure 7: 80% Consensus tree for bootstrapping hierarchical cluster on the samples, kidney data

## 5.2 An Affymetrix experiment

In this section we consider `abf1`, a 18 Affymetrix array experiment. Data is from 3 mouse strains, three individuals per strain and two arrays per individual. The experimental design is shown in Figure 8.

**Read Data** To reduce the size of *R/maanova* package, we do not include .CEL files,

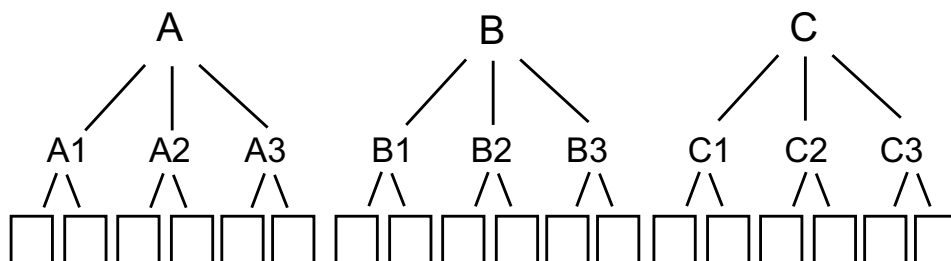


Figure 8: 18-array Affymetrix experiment design

thus you can not try following code. However this is the typical example to read the data from *affy* package. Suppose you have .CEL files and properly prepared TAB delimited text file under the working directory.

```
R> library(affy)
R> beforeRma <- ReadAffy()
R> rmaData <- rma(beforeRma)
R> datafile <- exprs(rmaData)
R> data1 <- read.madata(datafile=datafile,designfile="design.txt")
```

You want to save the output of `exprs()` to avoid reading the .CEL file again.

Another datafile that *R/maanova* can read is a TAB delimited text file.

```
R> data <- data.frame(probeid=row.names(datafile),datafile)
R> write.table(data,"data.dat",sep="\t",row=F,quote=F)
R> data2 <- read.madata(datafile="data.dat", designfile="design.txt",
  probeid=1, intensity=2)
```

**Design File** You must provide `designfile` information according to the order in which *affy* package reads the .CEL files, and save it as a TAB delimited text file. If you prepare `designfile` in advance, and want to read the .CEL file according to the order of `designfile`, you need to specify them at the filename in `ReadAffy()`. Here is the an example script that one can prepare `designfile`.

```
R> design.table <- data.frame(Array=row.names(pData(beforeRma)))
R> write.table(design.table,"design.txt",sep="\t",row=F,quote=F)
```

You can find 'design.txt' file, containing 'Array' field under the working directory. Open your favorite spread sheet and complete the rest of information. In this example, 'Strain' and 'Sample' information are provided. Table 1 shows the `designfile` of `abf1`. Note that `abf1` has `Sample` column. `Sample` representing biological replicate is not a required field in one color array, however it should be included in the `designfile`, when experiment has technical replicates.

Table 1: Design file for `abf1`

Array	Strain	Sample
E111.CEL	AJ	1
E112.CEL	AJ	1
E113.CEL	AJ	2
E114.CEL	AJ	2
E115.CEL	B6	3
E116.CEL	B6	3
...	...	...

Another option to prepare `designfile` is making a `data.frame` (or `matrix`) *R* object, and read it from *R* directly. This option is useful when the number of array is small. Again, make it sure that the information is provided by the order that `affy` reads the file.

```
R> design.table <- data.frame(Array=row.names(pData(beforeRma)));
R> Strain <- rep(c('AJ', 'B6', 'B6xAJ'), each=6)
R> Sample <- rep(c(1:9), each=2)
R> designfile <- cbind(design.table, Strain, Sample)
R> abf1 <- read.madata(datafile, designfile=designfile)
```

**Model fitting** This data has technical replicates, and thus `Sample` is treated as random term. `residplot()` shows residual plot.

```
R> data(abf1)
R> fit.full.mix <- fitmaanova(abf1, formula = ~Strain+Sample,
  random = ~Sample)
R> residplot(abf1, fit.full.mix)
```

**Test statistics** To test Strain effect in the model:

```
R> fttest.all = matest(abf1, fit.full.mix, test.method=c(1,1),
  shuffle.method="sample", term="Strain", n.perm= 100)
```

If you are interested in the pairwise comparison, you can specify the contrast matrix. Suppose you want to compare strain 1 vs. 2 and strain 1 vs. 3. Then the corresponding contrast matrix should have 2 rows (number of test) and three column (number of levels of term). You may want to use `PairContrast()` to get all possible pairwise comparisons.

```
R> C = matrix(c(1,-1,0,1,0,-1), ncol=3, byrow=T)
R> C.all = PairContrast(3)
R> ftest.pair = matest(abf1, fit.full.mix, Contrast = C,
  term="Strain", n.perm=100)
```

To derive FDR adjust p-value :

```
R> ftest.all = adjPval(ftest.all, method = 'jsFDR')
```

**Summarize the result** To get summary table, volcano plot and cluster :

```
R> summarytable(ftest.all)
R> summarytable(ftest.all, method=c("Pvalperm"), test =c("Fs"),
  whichTest=c("Fs.Pvalperm"), threshold=.1, outfile='shrotssummarytable.csv')

R> volcano(ftest.all)
R> idx.fix = volcano(ftest.pair)

R> cluster.kmean <- macluster(fit.full.mix, term="Strain",
  idx.gene=idx.fix$comparison1$idx.Fs,what="gene", method="kmean",
  kmean.ngroups=5, n.perm=100)
R> con.kmean <- consensus(cluster.kmean, 0.7)
R> con.kmean$groupname
R> cluster.hc <- macluster(fit.full.mix, term="Strain",
  idx.gene=idx.fix$comparison1$idx.Fs,what="sample", method="hc", n.perm=100)
R> con.hc <- consensus(cluster.hc)
```

Note that `ftest.pair` contains two comparisons, thus `volcano(ftest.pair)` generates two plots.

## References

- Y. Benjamin and Y. Hochberg. The control of the false discovery rate in multiple testing under dependency. *Ann. Stat.*, 29:1665–8, 2001.
- G.A. Churchill. Fundamentals of experimental design for cDNA microarrays. *Nat Genet*, 32(Suppl 2):490–5, 2002.
- XQ Cui and G.A. Churchill. Statistical Tests for Differential Expression in cDNA Microarray Experiments. *Genome Biology*, 4:201, 2003.

- XQ Cui, M.K. Kerr, and G.A. Churchill. Transformation for cDNA Microarray Data. *Statistical Applications in Genetics and Molecular Biology*, 2(1, article 4), 2003.
- XQ Cui, G. Hwang, J. Qiu, N. Blades, and G.A. Churchill. Improved Statistical Tests for Differential Gene Expression by Shrinking Variance Components. *Biostatistics*, 6: 59–75, 2005.
- J. Felsenstein. Confidence limits on phylogenies - An approach using the bootstrap. *Evolution*, 39:783, 1985.
- SAS institute Inc. *SAS/STAT User's Guide*. SAS institute Inc., 1999.
- M.K. Kerr and G.A. Churchill. Experimental design for gene expression microarrays. *Biostatistics*, 2:183, 2001a.
- M.K. Kerr and G.A. Churchill. Bootstrapping cluster analysis: Assessing the reliability of conclusion from microarray experiments. *PNAS*, 98:8961, 2001b.
- M.K. Kerr, M. Martin, and G.A. Churchill. Analysis of variance for gene expression microarray data. *J Computational Biology*, 7:819, 2000.
- M.K. Kerr, C.A. Afshari, L. Bennett, P. Bushel, J. Martinez, N. Walker, and G.A. Churchill. Statistical analysis of a gene expression microarray experiment with replication. *Statistica Sinica*, 12:203, 2001.
- M.K. Kerr, E.H. Leiter, L. Picard, and G.A. Churchill. *Computational and Statistical Approaches to Genomics*, chapter Sources of Variation in Microarray Experiments. Kluwer Academic Publishers, 2002.
- R.C. Littell, G.A. Milliken, W.W. Stroup, and R.D. Wolfinger. *SAS system for mixed models*. SAS institute Inc., 1996.
- T Margush and F.R. McMorris. Consensus n-trees. *Bulletin of Mathematical Biology*, 43:239, 1981.
- R.A. McLean, W.L. Sanders, and W.W. Stroup. A Unified Approach to Mixed Linear Models. *The American Statistician*, 45:54–64, 1991.
- C.C. Prichard, L. Hsu, J. Delrow, and P.S. Nelson. Project normal: defining normal variance in mouse gene expression. *PNAS*, 98:13266, 2001.
- S.R. Searle, G. Casella, and C.E. McCulloch. *Variance Components*. John Wiley and sons, Inc., 1992.
- V. Witkovsky. MATLAB algorithm mixed.m for solving Henderson's mixed model equations. *Technical Report, Inst. of Measurement Science, Slovak Academy of Science*, 2001.

R.D. Wolfinger, G. Gibson, E.D. Wolfinger, L. Bennett, H. Hamadeh, P. Bushel, C. Ashfari, and R.S. Paules. Assessing gene significance from cDNA microarray expression data via mixed models. *J Comp Biol*, 8:625, 2001.

Benjamin Y. and Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J Royal Statistical Society, Series B*(57):289, 1995.

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:e15, 2002.



## APPENDIX

### A Running R/maanova on computer cluster

Due to the intensive computation needed by the permutation test (especially for mixed effects models), parallel computing is implemented in R/maanova. The permutation tests can run on multiple computer nodes at the same time and the total computational time will be greatly reduced.

Here I provide some tips for system and software configurations for running R/maanova on clusters. My system is a 32-node beowulf cluster running Redhat linux 8.0 and R 1.8.0. You might encounter problems on different system following the steps.

#### A.1 System requirements

The computer cluster need to be Unix/Linux cluster with LAM/MPI installed. LAM/MPI can be download from <http://www.lam-mpi.org>. R and R/maanova need to be installed on all nodes (or they can be exported from a network file system).

The following R packages are required (you can obtain them from <http://cran.r-project.org>):

- SNOW (Simple Network Of Workstations). Obtain this one from <http://www.stat.uiowa.edu/~luke/R/cluster/cluster.html>.
- Rmpi
- serialize (this is required by Rmpi)

#### A.2 System setup

You might need to ask your system administrator to do this. But basically, you need to install LAM/MPI software in your system with correct configuration. Here are some tips for using rsh as remote shell program. If you want to use ssh, things will be slightly different. I assume LAM/MPI has been installed on your system.

1. First thing to do is to check if you can “rsh” to other machine without being asked for a password. If it does, add a file called “.rhosts” in your home directory. Each line of the file should be machine names and your user name. For a 2-node system, that file could look like (assume my account name is hao):

```
node1 hao
node2 hao
```

Then set the mode of this file to be 644 by doing “chmod 644 .rhosts”. After setting it up, do “rsh node1”. If it still asks for password, contact your system administrator.

2. Configure LAM to use rsh as the remote shell program by adding the following lines in your startup script (.bashrc, .cshrc, etc.)

```
LAMRSH="rsh"
```

```
export LAMRSH
```

(I found it's okay if you don't do this!)

3. Start LAM/MPI by typing "lamboot -v". You should see something like:

```
LAM 6.5.6/MPI 2 C++/ROMIO - University of Notre Dame
```

```
Executing hboot on n0 (node1 - 1 CPU)...
```

```
Executing hboot on n1 (node2 - 1 CPU)...
```

```
topology done
```

If you see some error messages, contact your system administrator.

### A.3 Install and test clusters in R

Here are the steps to install and use clusters in R.

1. Install all required package in R.
2. Start R, load in snow by doing "library(snow)".
3. To test if cluster is working, do the following:

```
R> library(snow)
```

```
R> cl <- makeCluster(2) # to make a cluster with two nodes
```

```
R> # to see if your cluster is correct
```

```
R> clusterCall(cl, function() Sys.info()[c("nodename", "machine")])
```

You're supposed to see the node names and machine types like what I got here:

```
[[1]]
```

```
nodename          machine
```

```
"node1"           "i686"
```

```
[[2]]
```

```
nodename          machine
```

```
"node2"           "i686"
```

4. Check the random number generations by running

```
clusterCall(cl, runif, 3)
```

If there is correlation problem, you must install rsprng.

That's it! Now you can start to enjoy the parallel computing.

## B Basic Algorithms

I will briefly introduce the core algorithms implemented in `R/maanova`. For the details please read related references.

### B.1 ANOVA model for microarray experiment

#### B.1.1 Model

The ANOVA model for microarray experiment was first proposed in Kerr *et al* (2000). Wolfinger *et al* (2001) proposed the 2-stage ANOVA model, which is the one used in `R/maanova`.

Basically an ANOVA model for microarray experiment can be specified in two stages. The first stage is the normalization model

$$y_{ijkgr} = \mu + A_i + D_j + AD_{ij} + r_{ijkgr} \quad (1)$$

The term  $\mu$  captures the overall mean. The rest of the terms capture the overall effects due to arrays, dyes, and labeling reactions. The residual of the first stage will be used as the inputs for the second stage.

The second stage models gene-specific effects:

$$r_{ijkgr} = G + AG_i + DG_j + VG_k + \epsilon_{ijkgr} \quad (2)$$

Here  $G$  captures the average effect of the gene.  $AG$  captures the array by gene variation.  $DG$  captures the dye by gene variation. For one dye system there will be no  $DG$  effect in the model.  $VG$  captures the effects for the experimental varieties. In the multiple factor experiment such like Paigen's 28-array experiment,  $VG$  should be further decoupled into several factors.

#### B.1.2 Fixed model versus Mixed model

The fixed effect model assumes independence among all observations and only one source of random variation. Although it is applicable to many microarray experiments, the fixed effects model does not allow for multiple sources of variation nor does it account for correlation among the observations that arise as a consequence of different layers of variation.

The mixed model treats some of the factors in an experimental design as random samples from a population. In other words, we assume that if the experiment were to be repeated, the same effects would not be exactly reproduced, but that similar effects would be drawn from a hypothetical population of effects. Therefore, we model these factors as sources of variance. Usually Array effect ( $AG$ ) should be treated as random factor in ANOVA model. If you have biological replicates, it should be treated as random as well.

## B.2 Hypothesis tests

For statistical test in the fixed ANOVA model, there is only one variance term and all factors in the model are tested against this variance. In mixed model ANOVA, there are multiple levels of variances (biological, array, residual, etc.). The test statistics need to be constructed based on the proper variations. I will skip the details here and the interested reader can read Searle *et al.*

### B.2.1 F test in matrix notation

Mathematically speaking, a Microarray ANOVA model for a gene can be expressed as the following mixed effect linear model (without losing generality):

$$y = X\beta + Zu + e \quad (3)$$

Here,  $\beta$  and  $u$  are vectors for fixed-effects and random-effects parameters.  $X$  and  $Z$  are design matrices. Note that fixed effects model is a special case for mixed effects model with  $Z$  and  $u$  empty. An assumption is that  $u$  and  $e$  are normally distributed with

$$E \begin{bmatrix} u \\ e \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad (4)$$

$$Var \begin{bmatrix} u \\ e \end{bmatrix} = \begin{bmatrix} G & 0 \\ 0 & R \end{bmatrix} \quad (5)$$

$G$  and  $R$  are unknown variance components and are estimated using restricted maximum likelihood (REML) method. The estimated  $\hat{\beta}$  and  $\hat{u}$  are called **best linear unbiased estimator** (BLUE) and **best linear unbiased predictor** (BLUP) respectively. Using estimated  $G$  and  $R$ , the variance-covariance matrix of  $\hat{\beta}$  and  $\hat{u}$  can be expressed as

$$\hat{C} = \begin{bmatrix} X'\hat{R}^{-1}X & X'\hat{R}^{-1}Z \\ Z'\hat{R}^{-1}X & Z'\hat{R}^{-1}Z + \hat{G}^{-1} \end{bmatrix}^{-} \quad (6)$$

where  $^{-}$  denotes generalized inverse. After getting  $\hat{C}$ , test statistics can be obtained through the following formula. For a given hypothesis

$$H : L \begin{bmatrix} \beta \\ u \end{bmatrix} = 0 \quad (7)$$

F-statistic is:

$$F = \frac{\begin{bmatrix} \hat{\beta} \\ \hat{u} \end{bmatrix}' L'(L\hat{C}L')^{-} L \begin{bmatrix} \hat{\beta} \\ \hat{u} \end{bmatrix}}{q} \quad (8)$$

Here  $q$  is the rank of  $L$  matrix.

Notice that for a fixed effect model where  $Z$  and  $u$  are empty, equation 8 will become

$$F = \frac{(L\hat{\beta})'(L(X'X)^{-1}L')^{-1}L\hat{\beta}}{q\sigma_e^2} \quad (9)$$

Here  $\sigma_e^2$  is the error variance.

F approximately follows F-distribution with  $q$  numerator degrees of freedom. The calculation of denominator degrees of freedom can be tricky and I will skip it here. If there is not enough degree of freedom, rely on tabulated P-values can be risky and a permutation test is recommended.

### B.2.2 Building $L$ matrix

In equation 8,  $L$  must be a matrix that  $L\beta$  is estimable. Estimability requires that the rows of  $L$  must be the linear combinations of the rows of  $X$ . In `R/maanova`, the program will automatically build  $L$  matrix for the term(s) to be tested. Normally in a mixed effects model, only the fixed effect terms can be tested so the  $u$  portion of  $L$  contains all 0s.

Building L matrix can be difficult. If the term to be tested is orthogonal to all other terms,  $L$  should contain all pairwise comparison for this term. That is, if the term has  $N$  levels, L should be a  $(N - 1)$  by  $N$  matrix. If the tested term is confounded with any other term, things will become complicated. `R/maanova` does the following to compute L matrix:

1. calculate a generalized inverse of  $X'X$  in such a way that the dependent columns in X for the tested term are set to zero.
2. calculate  $(X'X)^-(X'X)$ . The result matrix span the same linear space as  $X$  and it contains a lot of zeros.
3. Take the part of  $(X'X)^-(X'X)$  corresponding to the tested term, remove the rows with all zeros and the left part is  $L$ .

### B.2.3 Four flavors of F-tests

`R/maanova` offers four F statistics (called  $F_1$ ,  $F_2$ ,  $F_3$  and  $F_s$ ). Users have the option to turn on or off any of them. The difference among them is the calculation of  $\hat{C}$  matrix in equation 8. Or in another word, the way to estimate the variance components  $\hat{R}$  and  $\hat{G}$ .

Briefly speaking,  $F_1$  computes  $\hat{C}$  matrix based on the variance components of a single gene. It does not assume the common variance among the genes.  $F_3$  assumes the common variance among the genes. The  $\hat{C}$  matrix will be calculated based on the global variance (mean variance of all genes).  $F_2$  test is a hybrid of  $F_1$  and  $F_3$ . It uses a weighted combination of global and gene-specific variance to compute the  $\hat{C}$  matrix.  $F_s$  uses the James-Stein estimator for  $\hat{R}$  and  $\hat{G}$  and computes  $\hat{C}$  matrix. For details about four F-tests, please read Cui (2003b) and Cui (2003c).

### B.3 Data shuffling in permutation test

Choosing proper data shuffling method is crucial to permutation test. `R/maanova` offers two types of data shuffling, residual shuffling and sample shuffling.

Residual shuffling works only for fixed effects models. It is based on the assumption of homogeneous error variance. It does the following:

1. Fit a null hypothesis ANOVA model and get the residuals
2. Shuffle residuals globally and make a new data set
3. Compute test statistics on the new data set
4. Repeat step 2 and 3 for a certain times

Residual shuffling is incorrect for mixed effect models, where you have multiple random terms. For a mixed effects model, `R/maanova` will automatically choose sample shuffling, which does the following:

1. For the given tested term, check if it is nested within any random term. If not, go to step 4.
2. Choose the lowest random term (among the ones nesting with the tested term) as the base for shuffling. If there is only one random term nesting with the tested term, this random term will be the shuffle base.
3. Shuffle the sample names for the tested term in such a way that the same shuffle base will corresponds to the same sample name.
4. If there is no nesting, shuffle the sample names for the tested term freely, keeping all other terms unchanged. For multiple dye arrays, if Array and/or Dye effects are included in the model, the array structure should be preserved, e.g., the sample names on the same array should be shuffled together.
5. repeat step 3 or 4 for a certain times

Note that if the experiment size is small, the number of possible permutations will not be sufficient. In that case, users will have to rely on the tabulated values.

## C Frequently Asked Questions

1. Can R/*maanova* run on Windows XP?  
R/*maanova* runs under R. So if there is a version of R for your operating system, you can run R/*maanova*.
2. Can R/*maanova* handle missing data?  
No. R/*maanova* does not tolerant missing, zero or negative intensity data. Missing data will bring many problems. A gene with missing data will have a different experimental design from others. That will cause problem in permutation tests. So if you have any data missing, you must manually remove all data associated with that gene and all its replicates. We suggest you use non-background subtracted data as input.
3. Can R/*maanova* handle affymetric arrays?  
Yes. The newest version of R/*maanova* (0.97) works for N-dye system. But the data visualization and transformation functions are only working for 2-dye arrays at this time. So for affymetric data, I suggest you to use some other packages (such like *affy* package in BioConduction) to do the data transformation before loading into R/*maanova*.
4. Can R/*maanova* handle unbalanced design?  
Yes.
5. Can R/*maanova* output an ANOVA table for a model fitting?  
Not at this time. We used to output an ANOVA table in the older version of R/*maanova*. But as things getting complicated (with multiple factor design, mixed effects models, etc.), it becomes more and more difficult to build an ANOVA table. We think about it in the future.
6. What does R/*maanova* do with flagged genes?  
We did not do anything special to the flagged genes. They will be involved in data transformation and analysis. We just keep the flag in there so that if some of your significant genes were flagged, you need to put a question mark on it because the significance could be from a scratch on the slide.
7. Why I have problems reading in data?  
First of all data reading could behave differently in different operating system. Your data file need to be simple (ANSI) text file. Sometimes special characters (such like #, TAB, etc.) in the data file will cause problems. If you encounter that problem, remove the columns contains a lot of special characters (usually gene description) and try again. If you still have problem, contact the software author.
8. Why the estimated effects for a term don't sum to zero?  
When we build the design matrices we didn't put sum to zero constraints so that

the estimated effects for a term will not sum to zero. But the relative values of the estimates (which is the one we care about) will be correct.

9. Is there a graphical user interface for `R/maanova`?

Yes. We are developing a GUI for `R/maanova` using Java. It will be available in the near future.