

Exploring cDNA data

Nov 2004

Wolfgang Huber and Andreas Buess

The following exercise will show you some possibilities to load data from spotted cDNA microarrays into R, and to explore it using the statistical and visualization facilities of R and Bioconductor.

- 1.) **Preliminaries.** To go through this exercise, you need to have installed R \geq 2.0.0, the release 1.5 versions of the Bioconductor libraries Biobase, marray, multtest, limma, vsn, arrayMagic, and the library lymphoma, which contains the exercises and the lymphoma data set (see <http://lmpp.nih.gov/lymphoma/>).

```
> library("vsn")
> library("multtest")
> library("marray")
> library("limma")
> library("arrayMagic")
> library("lymphoma")
```

2.) Reading and exploring data files.

- a. First, you need to find the directory with the data on your harddisk. The path ends in `library/lymphoma/extdata`, and this is in the subdirectory where your R resides. On the command line, you can use the commands `dir()`, `getwd` and `setwd` to navigate around. In the GUI, you can use *File, Change dir* in the menu.
- b. Open the file `1c7b048rex.DAT` in a text editor. This is the typical file format for the results from the image analysis on a cDNA slide. Different image analysis programs use slightly different conventions and column headings, but you can always adapt the input function (see below) to your needs.
- c. Use the function `read.delim` to read the file into a *data frame* (that is a rectangular table of data) in R.

```
> x = read.delim("1c7b048rex.DAT")
```

- d. The table is too large to print it out as a whole, but we can find out about its size (with the function `dim`) and look at individual rows of the table.

```
> dim(x)
> colnames(x)
> x[1:6, ]
```

3.) Simple plots.

- a. Let us first look at the histogram of the values in the column `CH1I`, that is the channel 1 foreground intensity (see Fig. 1).

```
> hist(x$CH1I)
> hist(log2(x$CH1I), breaks = 100)
> hist(log2(x$CH1I), breaks = seq(5, 15, by = 0.25), col = "blue")
```

- b. Visualization of the spatial homogeneity of the hybridisations may help to assess their quality. The logarithm of the background of channel 1 is shown in Fig. 2.

```
> bg1 = spatialLayout(value = x$CH1B, row = x$ROW, col = x$COL,
+   block = x$GRID)
> plot(log(bg1), main = "Log Background of Channel 1")
```

We can also look at other transformations, like the logarithm or rank, to emphasize different kinds of inhomogeneities.

```
> plot(bg1)
> rankedbg1 = spatialLayout(value = rank(x$CH1B), row = x$ROW,
```

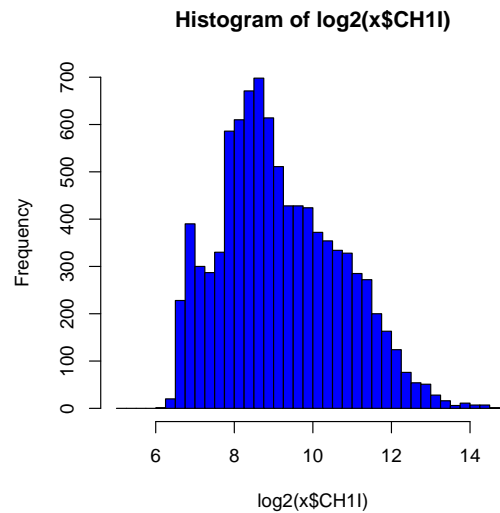


Figure 1:

```
+      col = x$COL, block = x$GRID)
> plot(rankedbg1)
Let us also have a look at the foreground of channel 1.
> fg1 = spatialLayout(value = x$CH1I, row = x$ROW, col = x$COL,
+      block = x$GRID)
> plot(log(fg1))
```

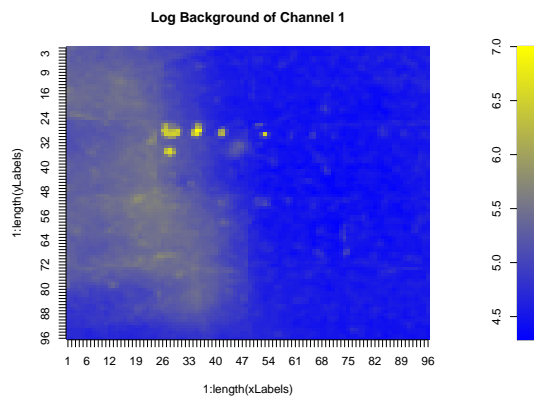


Figure 2:

c. Save one of the plots as PDF, and as Windows metafile. Copy and paste it into an MS-Office application.

4.) Calibration and variance stabilization.

a. Subtract the background intensities CH1B, CH2B from the foreground intensities CH1I, CH2I, and store the result in a 9216 x 2 matrix.

```
> y = cbind(x$CH1I - x$CH1B, x$CH2I - x$CH2B)
```

b. What does the function cbind do? Use the R online help to find out.

- c. Now we can use the function `vsn` to calibrate and transform the data, and plot the result (Fig. 3).

```
> ny = vsn(y)
> plot(exprs(ny), pch = ".")
```

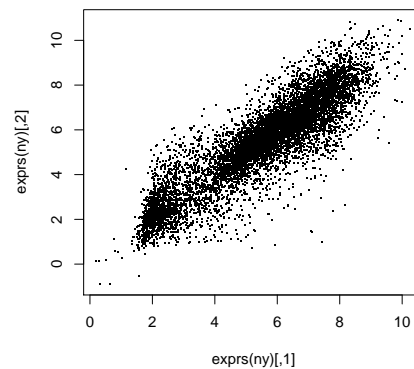


Figure 3:

- d. Have a look at the *vignette* — try the command `openVignette("vsn")`.

a. Reading a collection of files.

- b. The file `phenoData.txt` contains information on the samples that were hybridized onto the arrays. Look at it in a text editor. To load it into a `phenoData` object

```
> samples = read.phenoData("phenoData.txt", header = TRUE, as.is = TRUE)
> pData(samples)
```

	fileName	sampleid	tumortype	sex	slideNumber
1	lc7b047rex.DAT	CLL-13	CLL	m	1
2	lc7b048rex.DAT	CLL-13	CLL	m	2
3	lc7b069rex.DAT	CLL-52	CLL	f	3
4	lc7b070rex.DAT	CLL-39	CLL	f	4
5	lc7b019rex.DAT	DLCL-0032	DLCL	f	5
6	lc7b056rex.DAT	DLCL-0024	DLCL	m	6
7	lc7b057rex.DAT	DLCL-0029	DLCL	m	7
8	lc7b058rex.DAT	DLCL-0023	DLCL	<NA>	8

`phenoData` objects are where the Bioconductor stores information about samples, for example, treatment conditions in a cell line experiment or clinical or histopathological characteristics of tissue biopsies.

- c. Now we can load the whole set of 8 slides into the data object `a`.

```
> files = samples$fileName
> files
[1] "lc7b047rex.DAT" "lc7b048rex.DAT" "lc7b069rex.DAT" "lc7b070rex.DAT"
[5] "lc7b019rex.DAT" "lc7b056rex.DAT" "lc7b057rex.DAT" "lc7b058rex.DAT"
> theLayout = new("marrayLayout", maNgr = 4, maNgc = 4, maNsr = 24,
+   maNsc = 24, maNspots = 4 * 4 * 24 * 24)
> datdir = system.file("extdata", package = "lymphoma")
> a = read.marrayRaw(files, path = datdir, name.Gf = "CH1I", name.Gb = "CH1B",
+   name.Rf = "CH2I", name.Rb = "CH2B", layout = theLayout)
```

- d. ... and try out different normalization methods:

1. `vsN` (affine normalization and variance stabilization)
2. `maNorm` with global median location normalization
3. `maNorm` with loess for intensity- or A -dependent location normalization using the 'loess' smoother

```
> na1 = vsn(a)
> na2 = maNorm(a, norm = "median", echo = T)
```

- e. These commands take their time! You can save the results into a file with the `save` function, and later restore them with the `load` function. In MS-Windows, you can use the GUI for the latter.
- f. Now we want to extract the normalized log-ratios. The first line in the following code creates a three-dimensional array M with space for 9216 genes, 8 samples and 3 different normalization methods. `na1` is the result of `vsn`; the normalized intensities are accessed via the function `exprs`, and the log-ratios are obtained by subtracting the red intensities from the green ones. `na2` and `na3` are the output of `maarrayNorm`, and the log-ratios are obtained through the `slot` `maM` using the `accessor` `@`.

```
> odd = seq(1, 15, by = 2)
> even = seq(2, 16, by = 2)
> M = array(NA, dim = c(9216, 8, 2))
> M[, , 1] = exprs(na1)[, even] - exprs(na1)[, odd]
> M[, , 2] = na2@maM
```

- 5.) **Compare the results.** Look at scatterplots of the values of M from the same slide, calculated with different normalization methods. Do the values generally agree? How do they differ?

```
> plot(M[, 4, 1], M[, 4, 2], pch = ".", xlab = "vsn", ylab = "median")
```

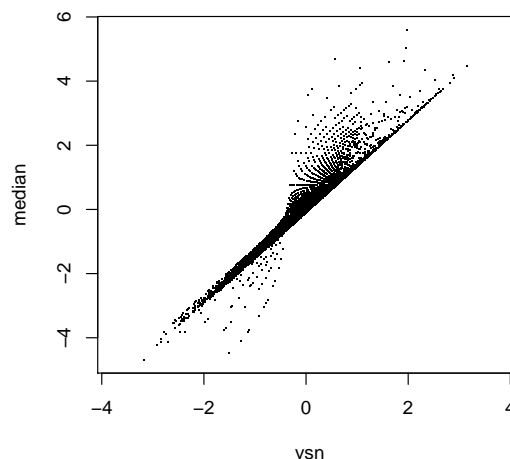


Figure 4:

- 6.) **Testing for differential transcription.** Now we are ready to calculate test statistics and to select genes. *Note:* The number of replicates (4 versus 4) that we are considering here is too small to derive significant conclusions about individual genes. The full data set contains many more chips. Here we restrict ourselves to a few of them in order to keep things simple for the purpose of this course.

- a. Look at the built-in function `t.test`, and at `mt.teststat` from the package `multtest`. Here, we use `mt.teststat` to calculate the t -test statistic for the comparison. The package `multtest` provides extensive functionality to calculate multiple-testing adjustments.

```
> classlabel = c(0, 0, 0, 0, 1, 1, 1, 1)
> tStat = mt.teststat(M[, , 1], classlabel)
> summary(tStat)
      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
-24.3400 -1.0330   0.1861   0.1832  1.3670  28.1000
> hist(tStat, breaks = 100, col = "#fb6090")
```

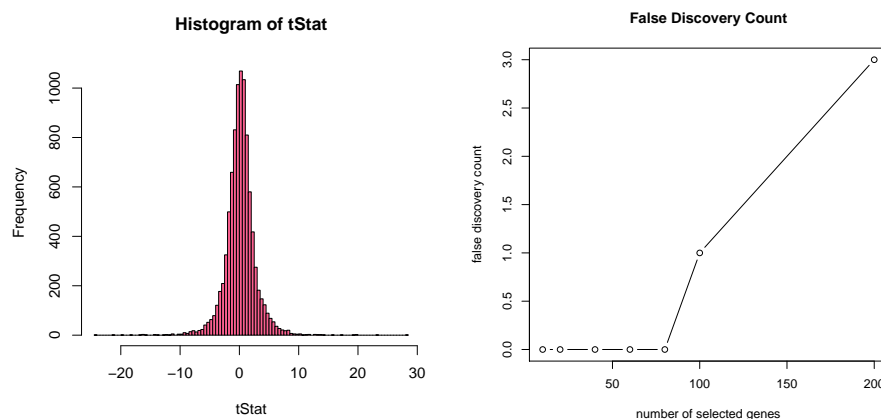


Figure 5: a) *left*: histogram of t -statistic, b) *right*: false discovery count.

- b. Similar to the FDR (false discovery rate) we estimate the significance of the most extreme t -values. The false discovery count is calculated for several lists of various length of the highest scored genes (Fig. 5 right).

```
> fc = fdc(M[, , 1], factor(classlabel), teststatfun = "rowttests")
> plot(fc$nrgenesel, fc$fdc, main = "False Discovery Count", xlab = "number of selected genes"
+      ylab = "false discovery count", type = "b")
```

- c. Now we load the spot (gene) description table

```
> spotDescr = read.delim("annotationData.txt")
```

and print the 5 genes with the lowest values of the t -statistic

```
> selection = order(tStat)[1:5]
> selection
```

```
[1] 4532 8076 6635 4586 739
```

as well as the 5 genes with the highest values of the t -statistic

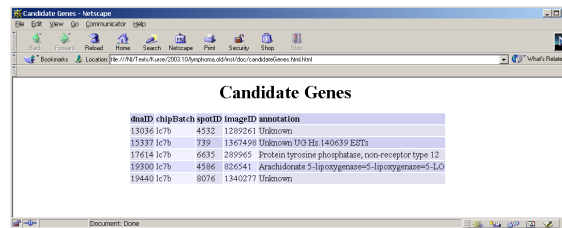
```
> selection = order(tStat, decreasing = TRUE)[1:5]
> selection
```

```
[1] 4323 4069 4331 2026 2143
```

In a following step we extract the annotation information for the 5 selected genes and generate a html-report (Fig. 6).

```
> spotIDs = x$SPOT[selection]
> geneAnno = spotDescr[spotDescr[, "spotID"] %in% spotIDs, ]
> write.htmltable(geneAnno, filename = "candidateGenes", title = "Candidate Genes")
```

7.) A simplified route to preprocessing and quality control.



The screenshot shows a web browser window titled "Candidate Genes - NetScape". The address bar shows a local file path. The main content area displays a table with the following data:

Read	chip	Batch	spotID	imageID	annotation
13036	ic7b		4512	1291241	Unknown
15337	ic7b	739	1367498	Unknown	U3 Hs 140639 ESTs
17614	ic7b	6635	289965		Protein tyrosine phosphatase, non-receptor type 12
19300	ic7b	4586	826541		Arachidonic 5-lipoxygenase=5-lipoxygenase=5-LO
19440	ic7b	8076	1340277		Unknown

Figure 6:

- a. The package arrayMagic can be used to automatically process the image files, to create R-objects and to generate quality diagnostics. See Figs. 7–9.

```
> res = processArrayData(slideDescriptionFile = "phenoData.txt",
+   loadPath = datdir, savePath = ".", fileNameColumn = "fileName",
+   slideNameColumn = "slideNumber", spotIdentifier = "SPOT",
+   type = "ScanAnalyze", subtractBackground = TRUE, normalisationMethod = "vsn")
> qR = qualityParameters(arrayDataObject = res$arrayDataObject,
+   exprSetRGObject = res$exprSetRGObject, spotIdentifier = "SPOT",
+   slideNameColumn = "slideNumber", resultFileName = "qualityResult.txt")
> qualityDiagnostics(exprSetRGObject = res$exprSetRGObject, arrayDataObject = res$arrayDataObject,
+   qualityParametersList = qR, plotOutput = "pdf", savePath = ".")
> visualiseHybridisations(arrayDataObject = res$arrayDataObject[,
+   c(1, 2, 6)], slideNameColumn = "slideNumber", mappingColumns = list(Block = "GRID",
+   Column = "COL", Row = "ROW"), type = "raw", savePath = ".",
+   plotOutput = "pdf")
> glratios = getExprSetLogRatio(res$exprSetRGObject)
> phenoD = phenoDataSlide(res$exprSetRGObject)
```

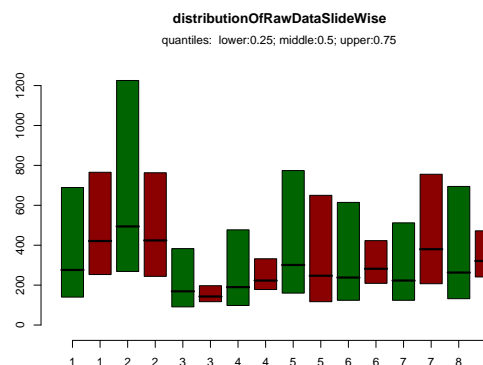


Figure 7: Interquartile ranges of raw data.

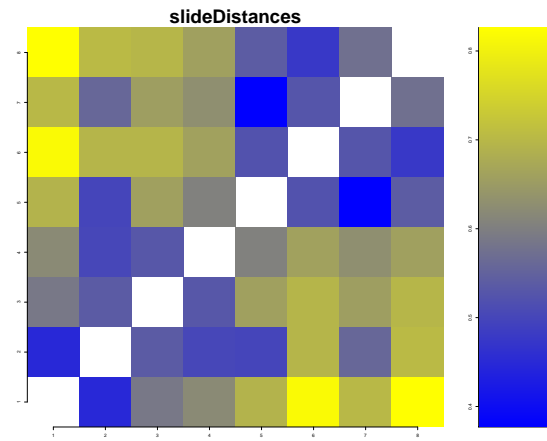


Figure 8: Matrix of distances between all pairs of slides.

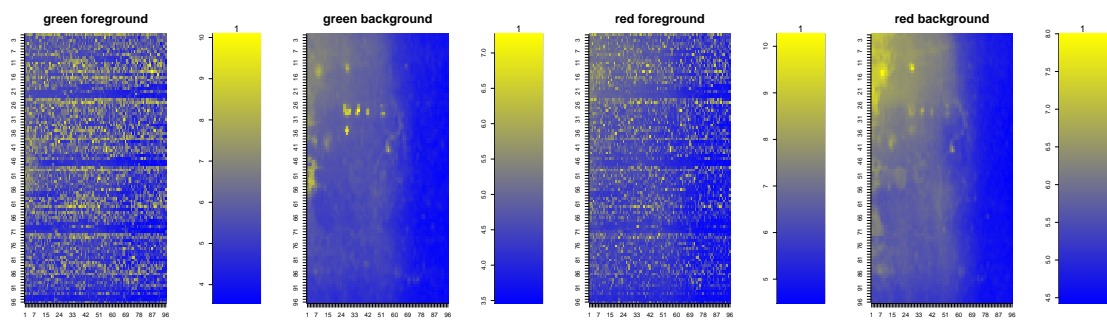


Figure 9: Spatial distribution of log raw intensities of one slide.