

Using the **a4** package

Willem Talloen, Tobias Verbeke

October 12, 2009

Contents

1 Preparation of the Data	2
1.1 ExpressionSet object	2
1.2 Some data manipulation	2
2 Unsupervised data exploration	3
3 Filtering	4
4 Detecting differential expression	5
4.1 T-test	5
4.2 Limma for comparing two groups	6
4.3 Limma for linear relations with a continuous variable	7
5 Class prediction	8
5.1 PAM	8
5.2 Random forest	8
5.3 Forward filtering with various classifiers	9
5.4 Penalized regression	11
5.5 Logistic regression	13
5.6 Receiver operating curve	15
6 Visualization of interesting genes	16
6.1 Plot the expression levels of one gene	16
6.2 Plot the expression levels of two genes versus each other	20
6.3 Plot expression line profiles of multiple genes/probesets across samples	21
6.4 Smoothscatter plots	23
6.5 Gene lists of log ratios	26
7 Pathway analysis	27
7.1 Minus log p	27
7.2 Gene set enrichment analysis	27
8 Software used	27

1 Preparation of the Data

First we load the package `a4` and the example real-life data set `ALL`.

```
R> library(a4)
R> data(ALL, package = "ALL")
```

For illustrative purposes, simulated data sets can also be very valuable (but not used here).

```
R> esSim <- simulateData(nEffectRows = 50, betweenClassDifference = 5,
  nNoEffectCols = 5, withinClassSd = 0.2)
```

1.1 ExpressionSet object

The data are assumed to be in an expressionSet object. Such an object structure combines different sources of information into a single structure, allowing easy data manipulation (e.g., subsetting, copying) and data modelling.

The `textttfeatureData` slot is typically not yet containing all relevant information about the genes. This interesting extra gene information can be added using `addGeneInfo`.

```
R> ALL <- addGeneInfo(ALL)
```

1.2 Some data manipulation

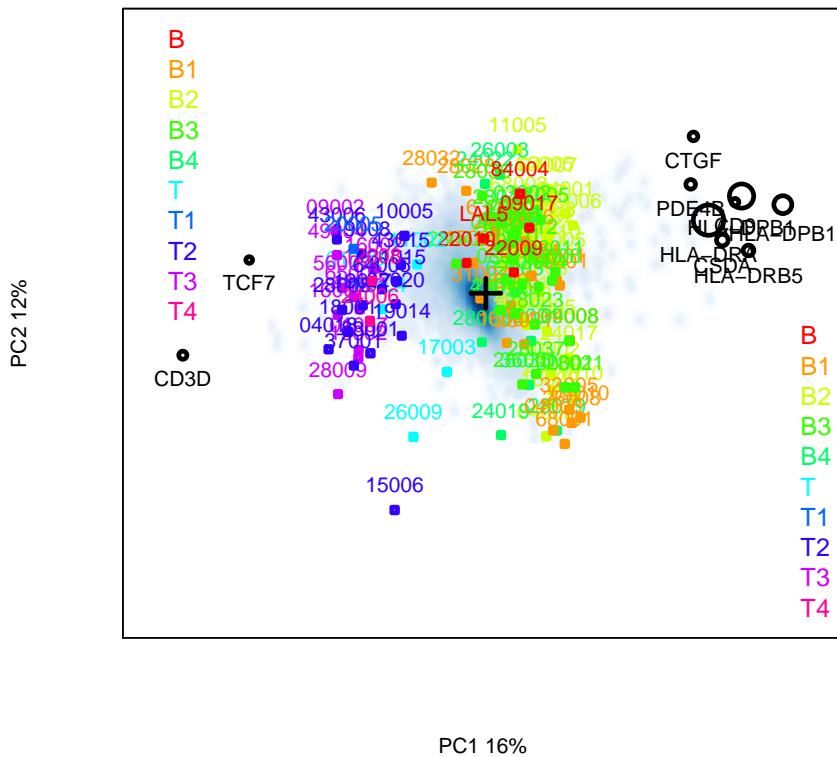
The `ALL` data consists out of samples obtained from two types of cells with very distinct expression profiles; B-cells and T-cells. To have a more subtle signal, gene expression will also be compared between the BCR/ABL and the NEG group within B-cells only. To this end, we create the expressionSet `bcrAb1OrNeg` containing only B-cells with BCR/ABL or NEG.

```
R> Bcell <- grep("^B", as.character(ALL$BT))
R> subsetType <- "BCR/ABL"
R> bcrAb1OrNegIdx <- which(as.character(ALL$mol) %in%
  c("NEG", subsetType))
R> bcrAb1OrNeg <- ALL[, intersect(Bcell, bcrAb1OrNegIdx)]
R> bcrAb1OrNeg$mol.biol <- factor(bcrAb1OrNeg$mol.biol)
```

2 Unsupervised data exploration

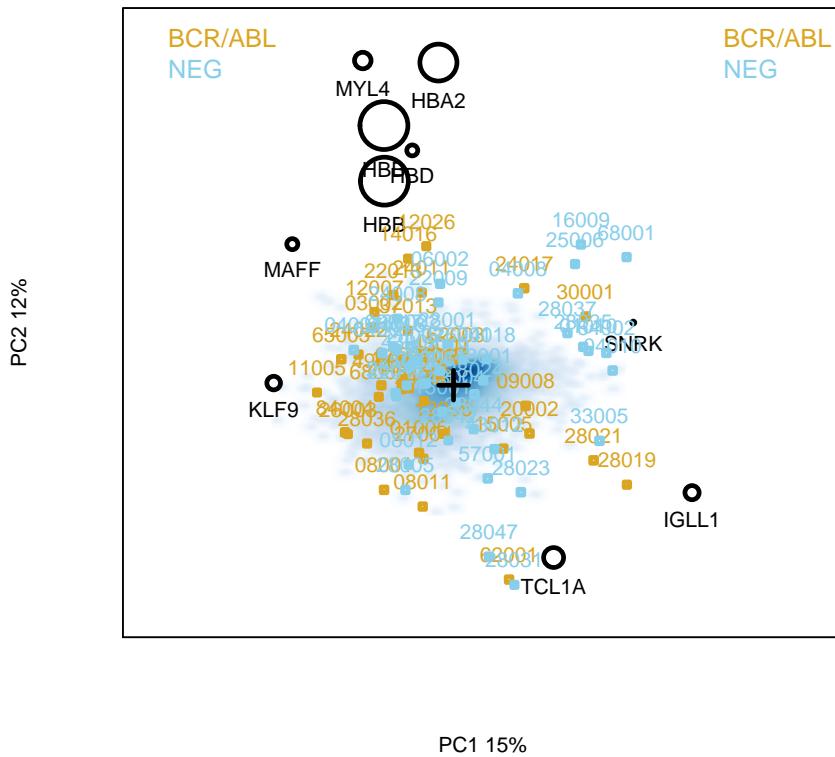
Spectral maps are very powerful techniques to get an unsupervised picture of how the data look like. A spectral map of the ALL data set shows that the B- and the T-subtypes cluster together along the x-axis (the first principal component). The plot also indicates which genes contribute in which way to this clustering. For example, the genes located in the same direction as the T-cell samples are higher expressed in these T-cells. Indeed, the two genes at the left (TCF7 and CD3D) are well known to be specifically expressed by T-cells (Wetering 1992, Krissansen 1986).

```
R> spectralMap(object = ALL, groups = "BT")
R> legend("bottomright", legend = levels(pData(ALL)$BT),
  text.col = a4palette(nlevels(pData(ALL)$BT)),
  bty = "n")
```



A spectral map of the `bcrAb1OrNeg` data subset does not show a clustering of BCR/ABL or NEG cells.

```
R> spectralMap(object = bcrAb1OrNeg, groups = "mol.biol",
  probe2gene = TRUE)
R> legend("topright", legend = levels(pData(bcrAb1OrNeg)$mol.biol),
  text.col = a4palette(nlevels(pData(bcrAb1OrNeg)$mol.biol)),
  bty = "n")
```



3 Filtering

The data can be filtered, for instance based on variance and intensity, in order to reduce the high-dimensionality.

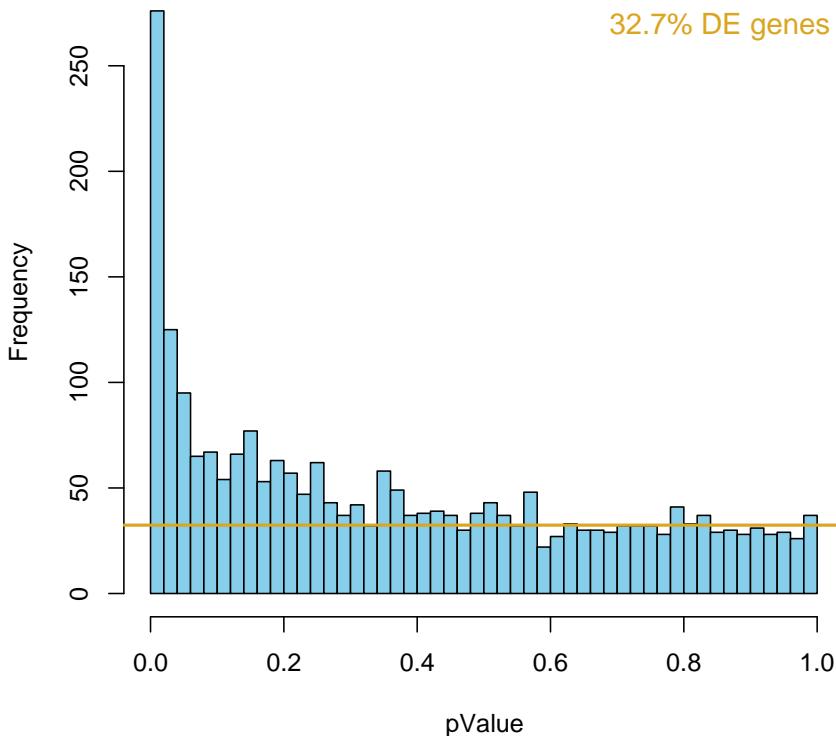
```
R> selBcrAb1OrNeg <- filterVarInt(object = bcrAb1OrNeg)
R> propSelGenes <- round((dim(selBcrAb1OrNeg)[1]/dim(bcrAb1OrNeg)[1]) *
  100, 1)
```

This filter selected 18.9 % of the genes (2391 of the in total 12625 genes).

4 Detecting differential expression

4.1 T-test

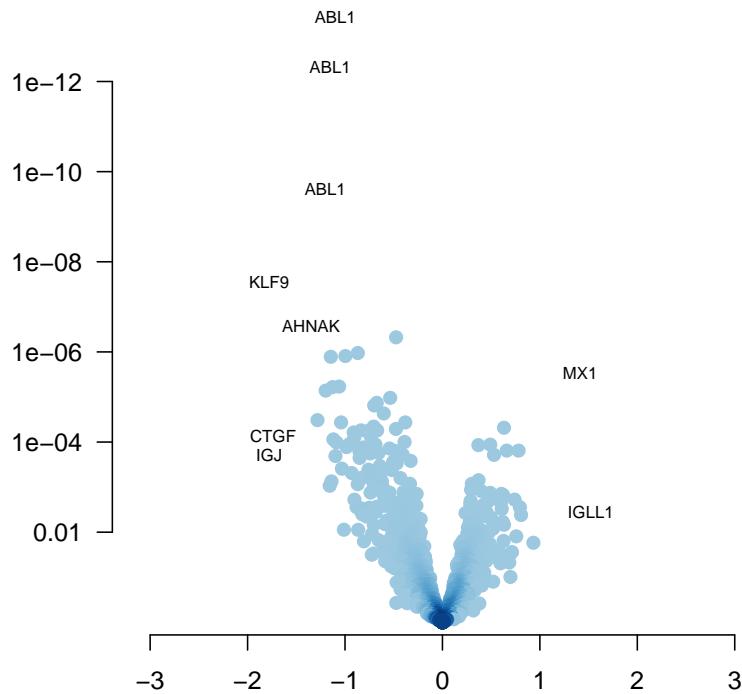
```
R> tTestResult <- tTest(selBcrAb1OrNeg, "mol.biol")  
R> histPvalue(tTestResult[, "p"], addLegend = TRUE)  
R> propDEgenesRes <- propDEgenes(tTestResult[,  
  "p"])
```



Using an ordinary t-test, there are 171 genes significant at a FDR of 10%. The proportion of genes that are truly differentially expressed is estimated to be around 32.7.

The toptable and the volcano plot show that three most significant probe sets all target **ABL1**. This makes sense as the main difference between BCR/ABL and NEG cells is a mutation in this particular ABL gene.

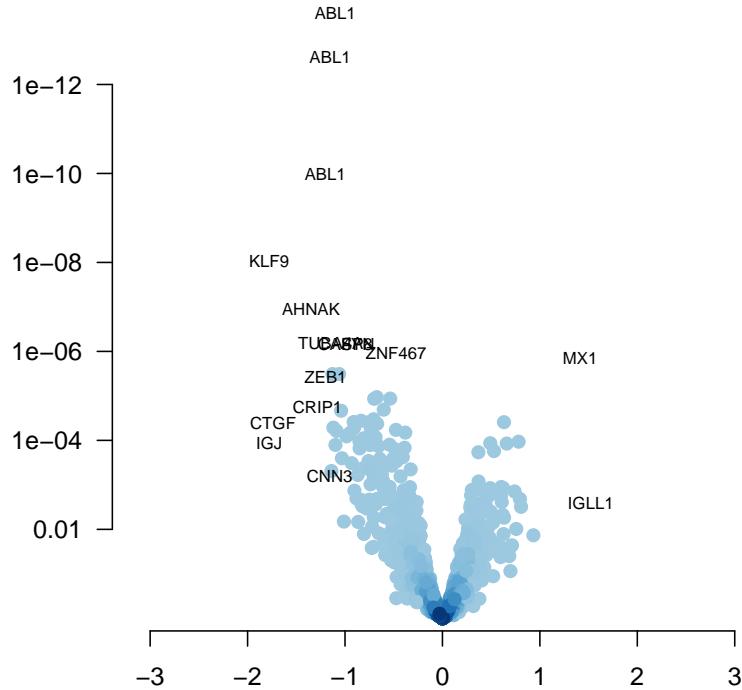
```
R> tabTTest <- topTable(tTestResult, n = 10)  
R> xtable(tabTTest, caption = "The top 5 features selected by an ordinary t-test.",  
  label = "tablassoClass")  
  
R> volcanoPlot(tTestResult, topPValues = 5, topLogRatios = 5)
```



4.2 Limma for comparing two groups

In this particular data set, the modified t-test using `limma2Groups` provides very similar results. This is because the sample size is relatively large.

```
R> limmaResult <- limma2Groups(selBcrAb10rNeg,
  "mol.biol")
R> volcanoPlot(limmaResult)
```



It is very useful to put lists of genes in annotated tables where the genes get hyperlinks to EntrezGene.

```
R> tabLimma <- topTable(limmaResult, n = 10)
```

Gene	logFC	AveExpr	P.Value	adj.P.Val	Description
ABL1	-1.10	9.20	0.00	0.00	c-abl oncogene 1, receptor tyrosine ki
ABL1	-1.15	9.00	0.00	0.00	c-abl oncogene 1, receptor tyrosine ki
ABL1	-1.20	7.90	0.00	0.00	c-abl oncogene 1, receptor tyrosine ki
KLF9	-1.78	8.62	0.00	0.00	Kruppel-like factor 9
AHNAK	-1.35	8.44	0.00	0.00	AHNAK nucleoprotein
TUBA4A	-1.15	9.23	0.00	0.00	tubulin, alpha 4a
FYN	-0.87	7.76	0.00	0.00	FYN oncogene related to SRC, FGR, YES
CASP8	-1.00	8.04	0.00	0.00	caspase 8, apoptosis-related cysteine
ZNF467	-0.48	7.14	0.00	0.00	zinc finger protein 467
MX1	1.41	6.73	0.00	0.00	myxovirus (influenza virus) resistance

4.3 Limma for linear relations with a continuous variable

Testing for (linear) relations of gene expression with a (continuous) variable is typically done using regression. A modified t-test approach improves the results by penalizing small slopes. The modified regressions can be applied using `limmaReg`.

5 Class prediction

There are many classification algorithms with profound conceptual and methodological differences. Given the differences between the methods, there's probably no single classification method that always works best, but that certain methods perform better depending on the characteristics of the data.

On the other hand, these methods are all designed for the same purpose, namely maximizing classification accuracy. They should consequently all pick up (the same) strong biological signal when present, resulting in similar outcomes.

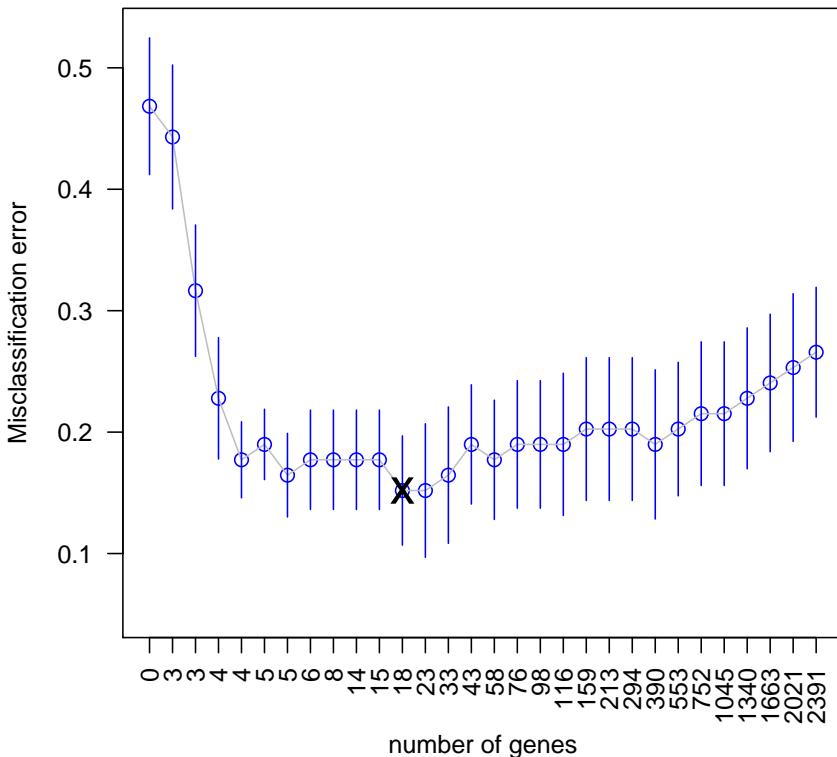
Personally, we like to apply four different approaches; PAM, RandomForest, forward filtering in combination with various classifiers, and LASSO.

All four methods have the property that they search for the smallest set of genes while having the highest classification accuracy. The underlying rationale and algorithm is very different between the four approaches, making their combined use potentially complementary.

5.1 PAM

PAM (Tibshirani 2002) applies univariate and dependent feature selection.

```
R> resultPam <- pamClass(selBcrAb1OrNeg, "mol.biol")
R> plot(resultPam)
R> featResultPam <- topTable(resultPam, n = 15)
R> xtable(head(featResultPam$listGenes), caption = "Top 5 features selected by PAM.")
```



5.2 Random forest

Random forest with variable importance filtering (Breiman 2001, Diaz-Uriarte 2006) applies multivariate and dependent feature selection. Be cautious when interpreting its outcome, as the obtained results are unstable and sometimes overoptimistic.

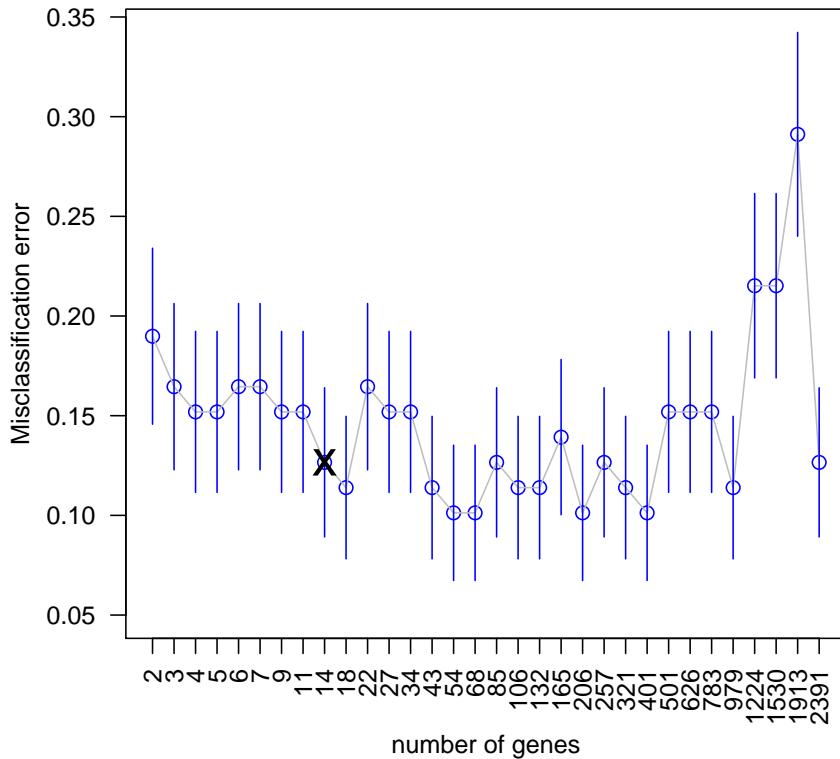
```

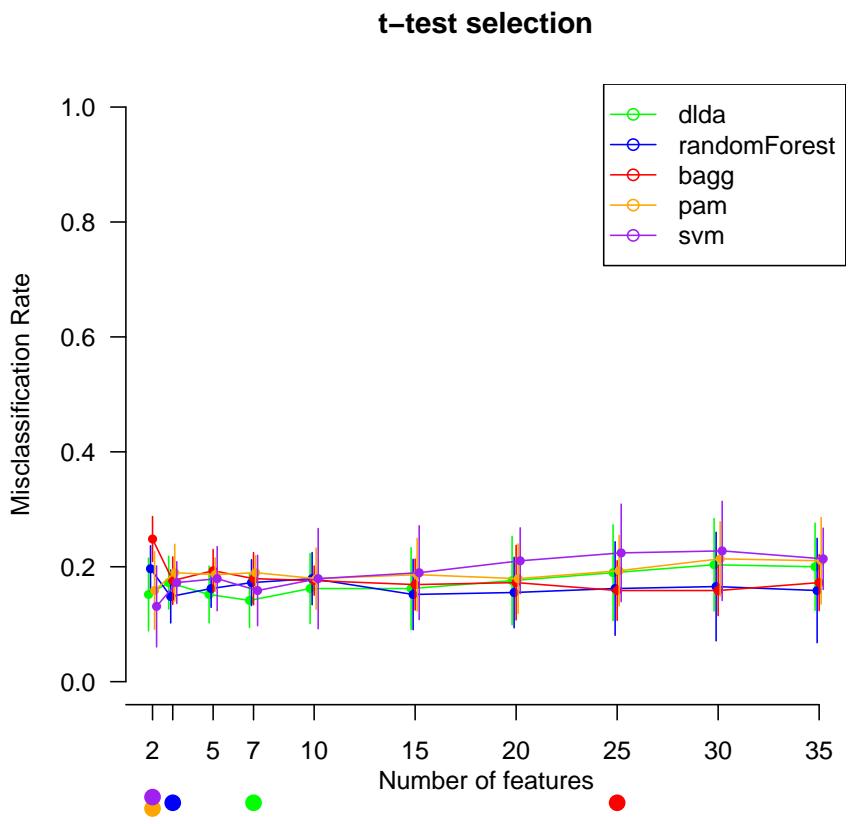
R> resultRF <- rfClass(selBcrAb10rNeg, "mol.biol")
R> plot(resultRF, which = 2)
R> featResultRF <- topTable(resultRF, n = 15)
R> xtable(head(featResultRF$topList), caption = "Features selected by Random Forest variable importance")

```

	GeneSymbol
1635_at	ABL1
1636_g_at	ABL1
2056_at	FGFR1
32116_at	TMC6
33774_at	CASP8
37027_at	AHNAK

Table 1: Features selected by Random Forest variable importance.



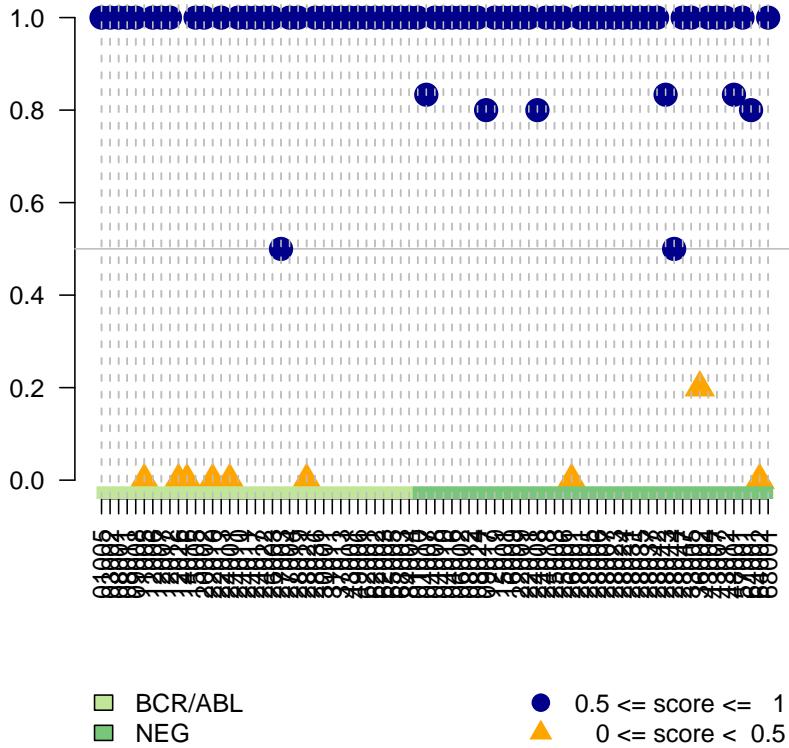


	nFeat_optim	mean_MCR	sd_MCR
dlda	7.00	0.14	0.05
randomForest	3.00	0.15	0.05
bagg	25.00	0.16	0.05
pam	2.00	0.16	0.07
svm	2.00	0.13	0.07

Table 2: Optimal number of genes per classification method together with the respective misclassification error rate (mean and standard deviation across all CV loops).

```
R> scoresPlot(nlcvTT, tech = "svm", nfeat = 2)
```

Scores Plot (svm, 2 features)

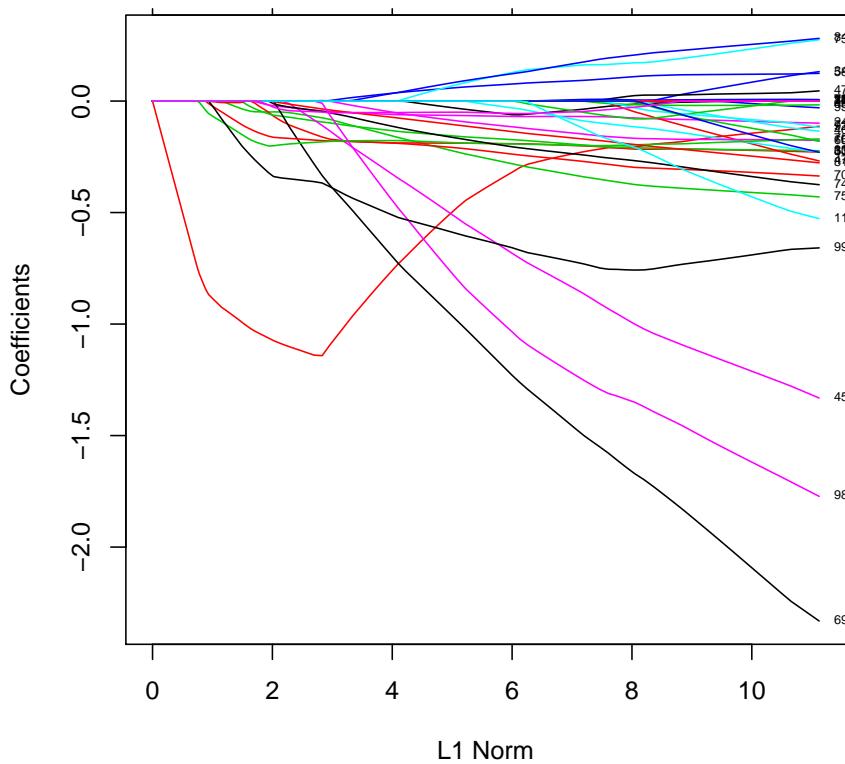


5.4 Penalized regression

LASSO (Tibshirani 2002) or elastic net (Zou 2005) apply multivariate and dependent feature selection.

```
R> resultLasso <- lassoClass(object = bcrAb1OrNeg, groups = "mol.biol")
R> plot(resultLasso, label = TRUE, main = "Lasso coefficients in relation to degree of penalization")
R> featResultLasso <- topTable(resultLasso, n = 15)
```

Lasso coefficients in relation to degree of penalization.



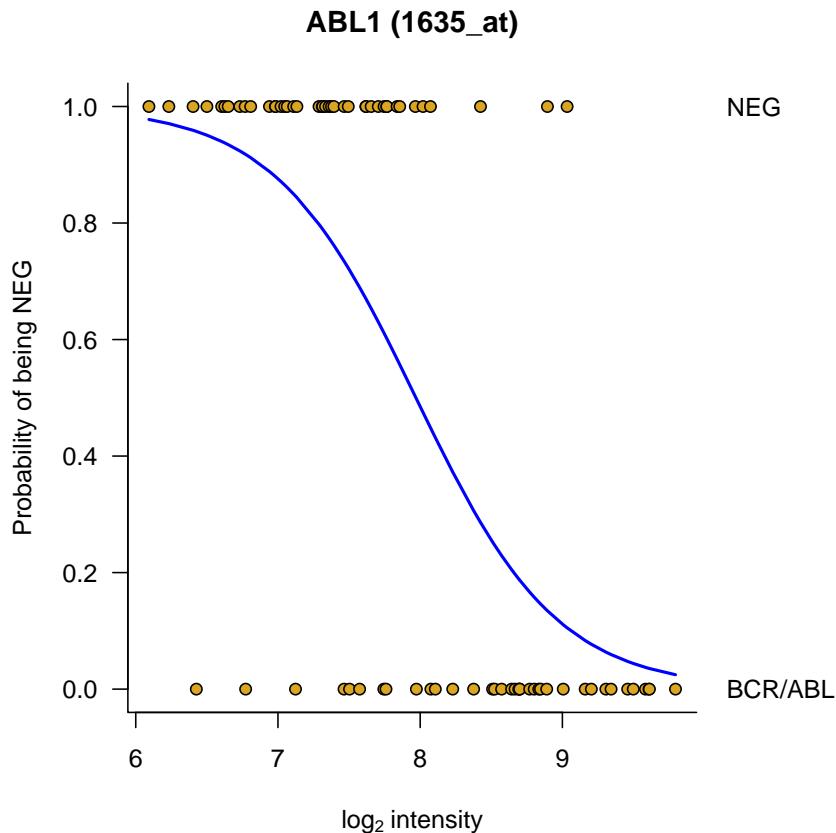
Gene	Coefficient
ITGA7	-2.33
ABL1	-1.77
TCL1B	-1.33
ZNF467	-0.66
RAB32	-0.53
YES1	-0.43
ANXA1	-0.38
ALDH1A1	-0.34
DSTN	0.28
F13A1	-0.28
PTDSS1	0.27
CHD3	-0.27
SERPINE2	-0.23
GAB1	-0.23
TLR1	-0.23

Table 3: Features selected by Lasso, ranked from largest to smallest penalized coefficient.

5.5 Logistic regression

Logistic regression is used for predicting the probability to belong to a certain class in binary classification problems.

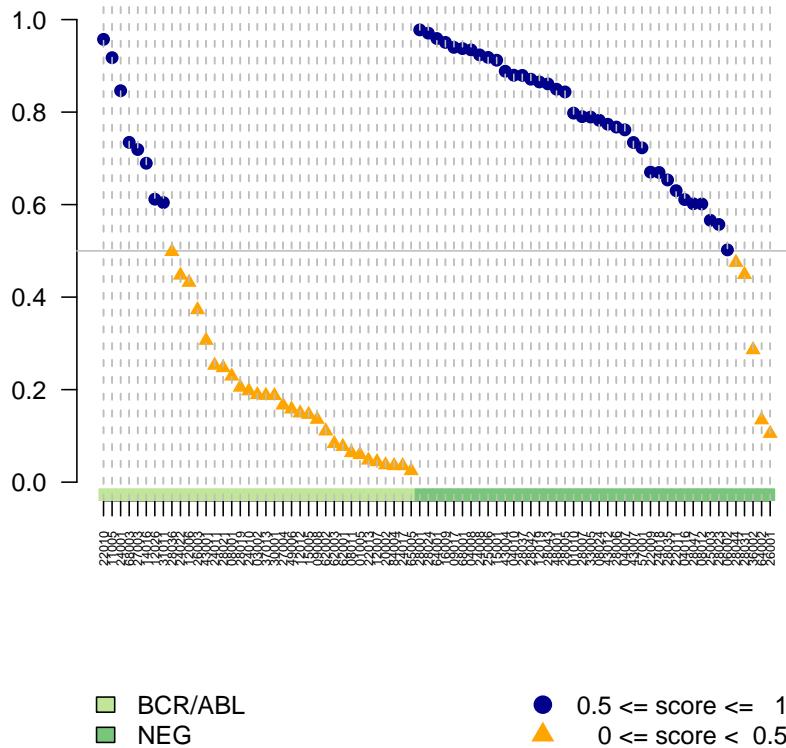
```
R> logRegRes <- logReg(geneSymbol = "ABL1", object = bcrAblOrNeg, groups = "mol.biol")
```



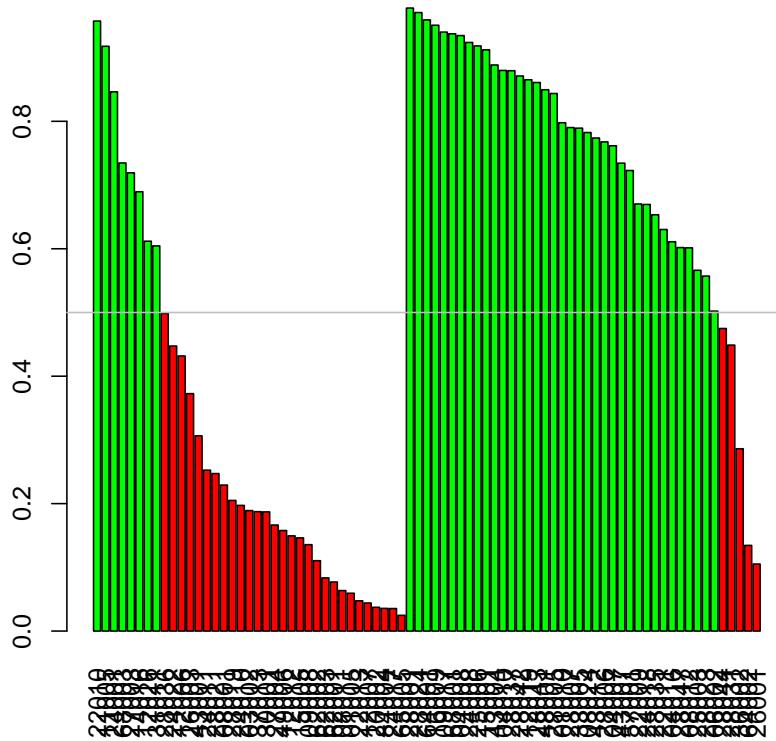
The obtained probabilities can be plotted with `ProbabilitiesPlot`. A horizontal line indicates the 50% threshold, and samples that have a higher probability than 50% are indicated with blue dots. Apparently, using the expression of the gene **ABL1**, quite a lot of samples predicted to with a high probability to be NEG, are indeed known to be NEG.

```
R> probabilitiesPlot(proportions = logRegRes$fit, classVar = logRegRes$y, sampleNames = rownames(1)
```

Probability of being NEG



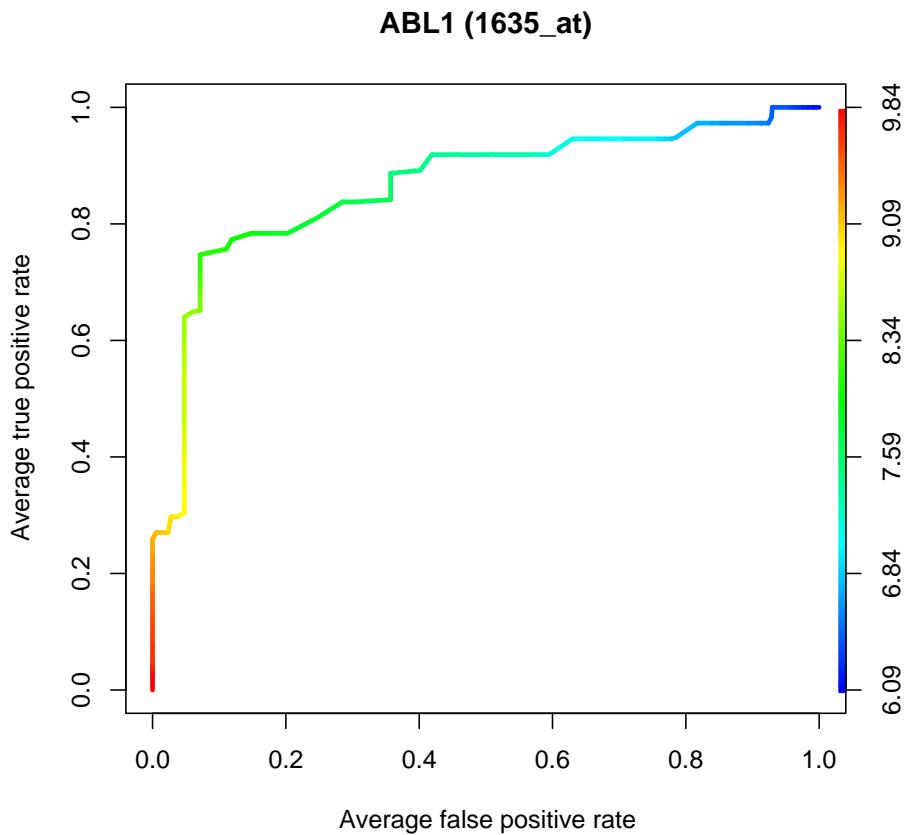
```
R> probabilitiesPlot(proportions = logRegRes$fit, classVar = logRegRes$y, barPlot = TRUE, sampleNa
```



5.6 Receiver operating curve

A ROC curve plots the fraction of true positives (TPR = true positive rate) versus the fraction of false positives (FPR = false positive rate) for a binary classifier when the discrimination threshold is varied. Equivalently, one can also plot sensitivity versus (1 - specificity).

```
R> ROCres <- ROCcurve(geneSymbol = "ABL1", object = bcrAb1OrNeg, groups = "mol.biol")
```



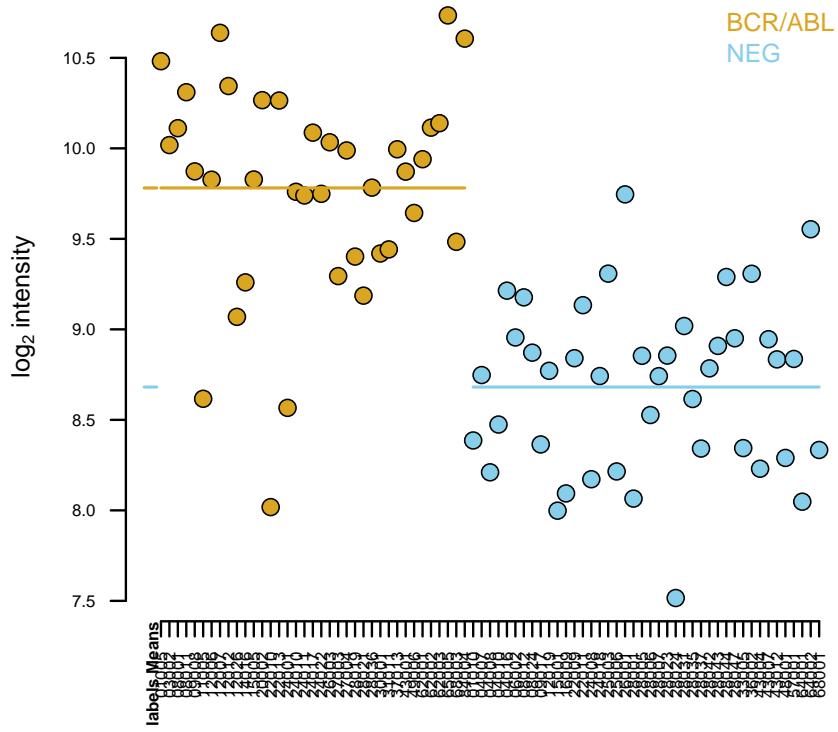
6 Visualization of interesting genes

6.1 Plot the expression levels of one gene

Some potentially interesting genes can be visualized using `plot1gene`. Here the most significant gene is plotted.

```
R> plot1gene(probesetId = rownames(tTestResult)[1], object = selBcrAbl0rNeg, groups = "mol.biol",
```

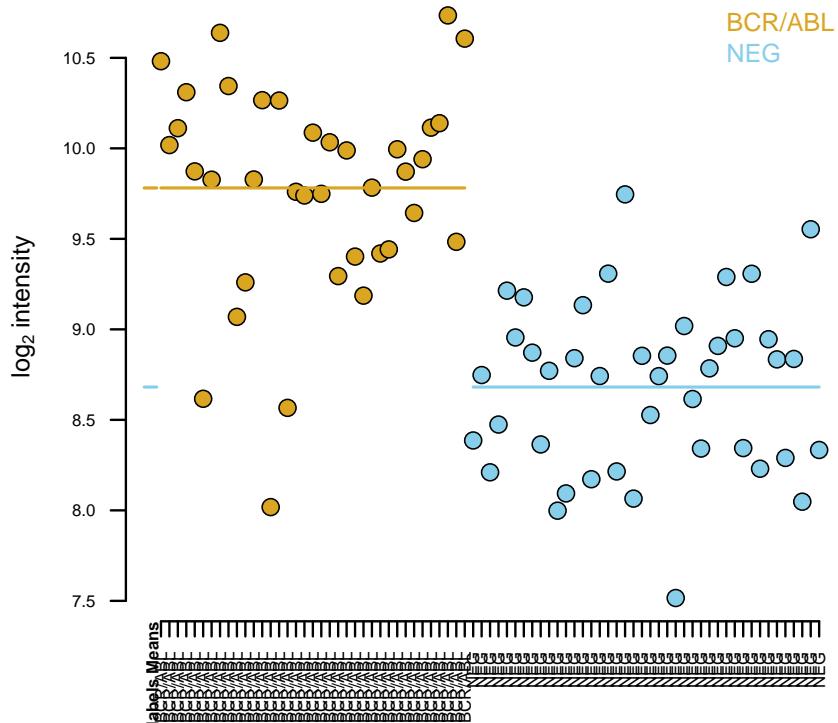
ABL1 (1636_g_at)



There are some variations possible on the default `plot1gene` function. For example, the labels of x-axis can be changed or omitted.

```
R> plot1gene(probesetId = rownames(tTestResult)[1], object = selBcrAb1OrNeg, groups = "mol.biol",
```

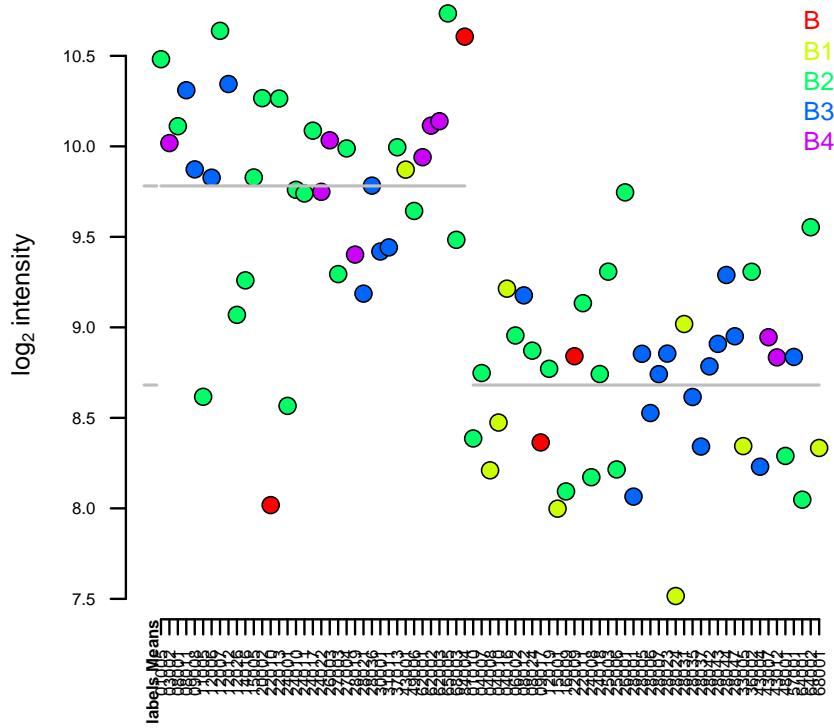
ABL1 (1636_g_at)



Another option is to color the samples by another categorical variable than used for ordering.

```
R> plot1gene(probesetId = rownames(tTestResult)[1], object = selBcrAb1OrNeg, groups = "mol.biol",
```

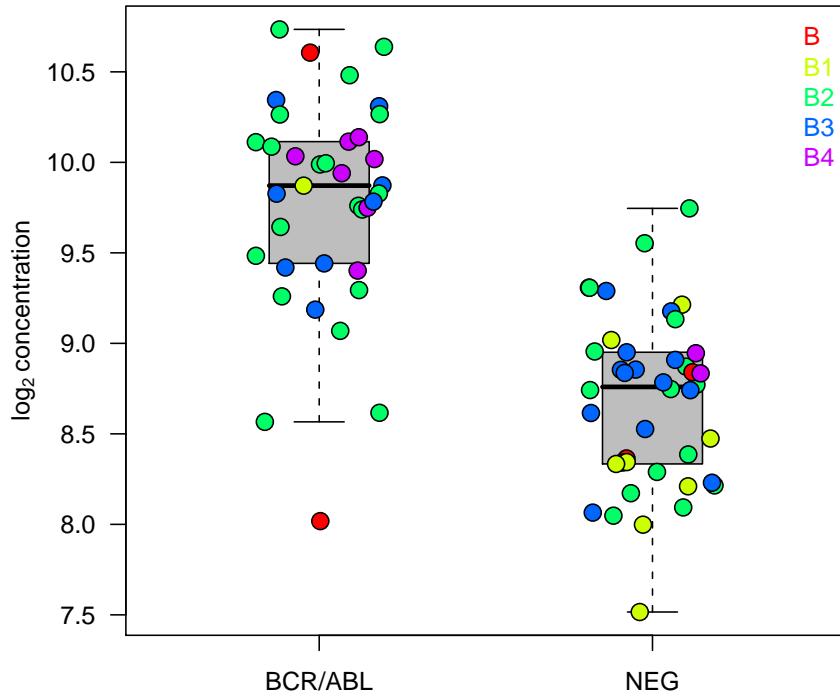
ABL1 (1636_g_at)



The above graphs plot one sample per tickmark in the x-axis. This is very useful to explore the data as one can directly identify interesting samples. If it is not interesting to know which sample has which expression level, one may want to plot in the x-axis not the samples but the groups of interest. It is possible to pass arguments to the boxplot function to customize the graph. For example the `boxwex` argument allows to reduce the width of the boxes in the plot.

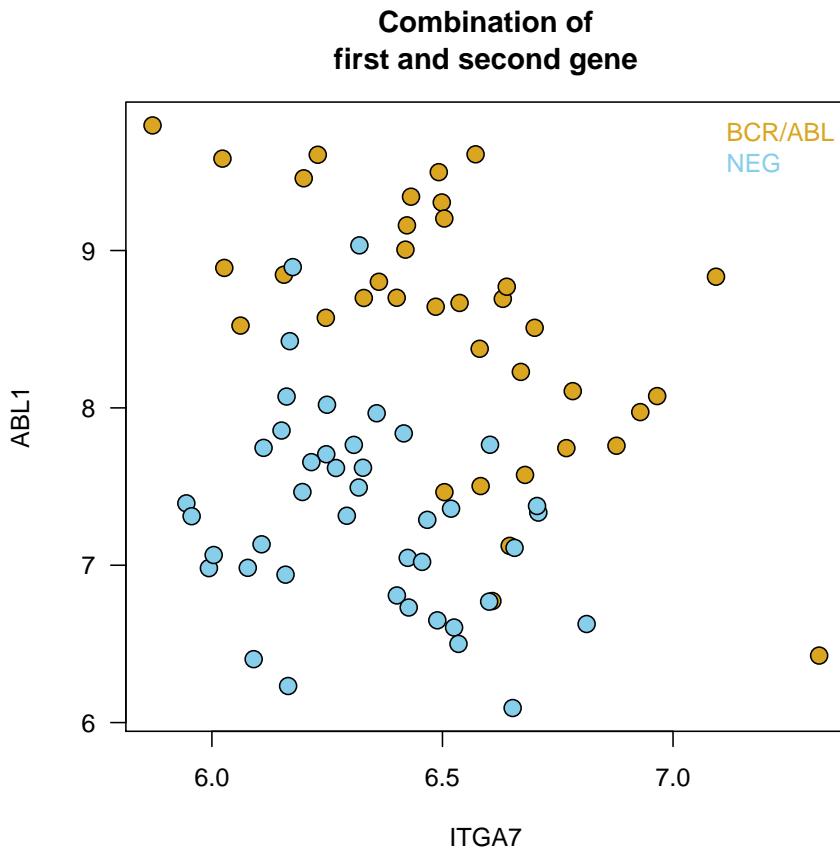
```
R> boxPlot(probesetId = rownames(tTestResult)[1], object = selBcrAb1OrNeg, boxwex = 0.3, groups =
```

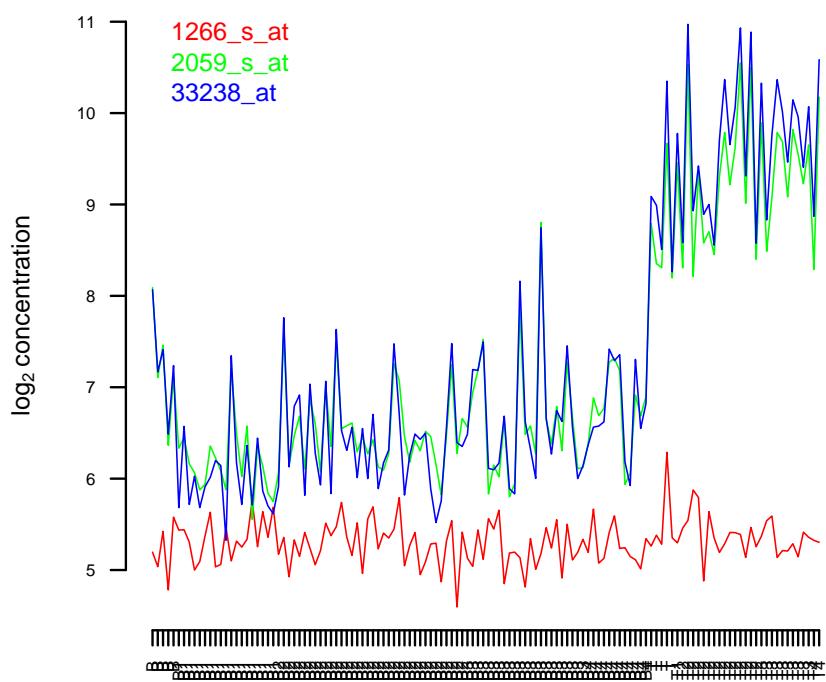
ABL1 (1636_g_at)



6.2 Plot the expression levels of two genes versus each other

```
R> plotCombination2genes(geneSymbol1 = featResultLasso$topList[1, 1], geneSymbol2 = featResultLasso$topList[2, 1])
```





6.4 Smoothscatter plots

It may be of interest to look at correlations between samples. As each dot represents a gene, there are typically many dots. It is therefore wise to color the dots in a density dependent way.

```
R> plotComb2Samples(ALL, "11002", "01003", xlab = "a T-cell", ylab = "another T-cell")
```

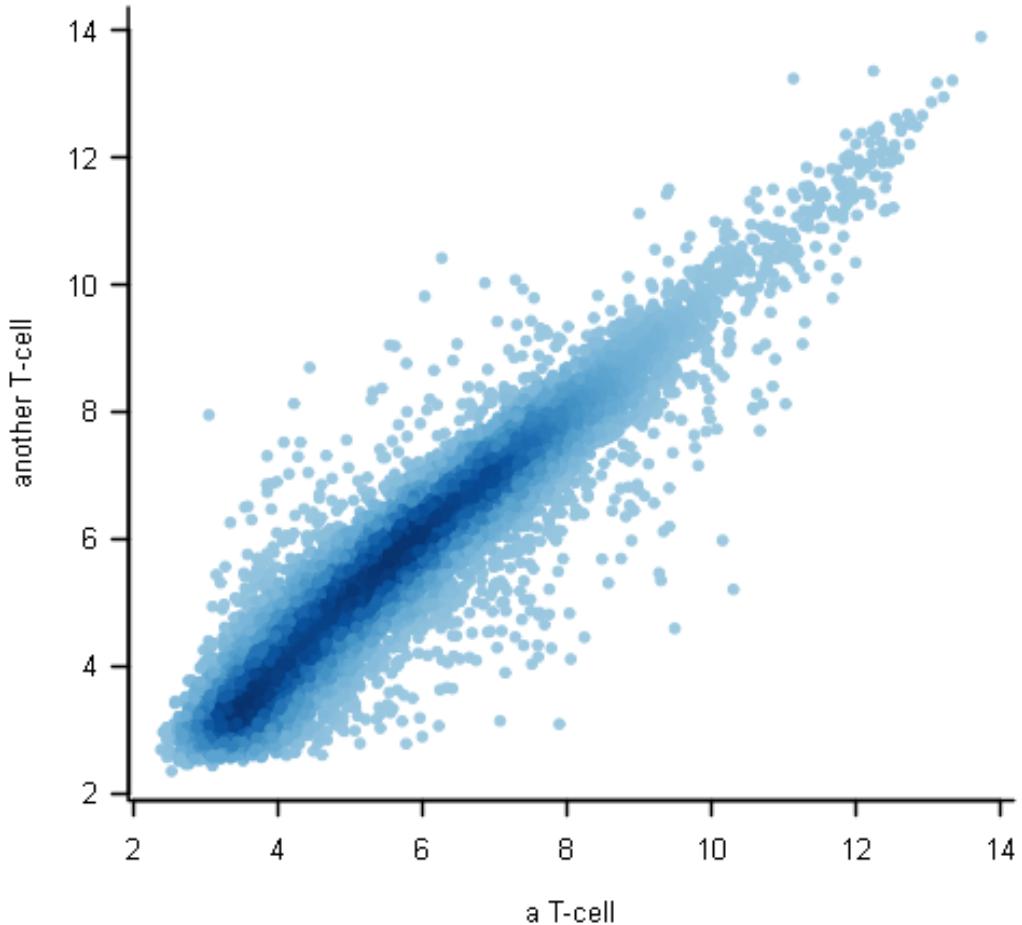


Figure 1: Correlations in gene expression profiles between two T-cell samples (samples 11002 and 01003).

If there are outlying genes, one can label them by their gene symbol by specifying the expression intervals (X- or Y- axis or both) that contain the genes to be highlighted using `trsholdX` and `trsholdY`.

```
R> plotComb2Samples(ALL, "84004", "01003", trsholdX = c(10, 12), trsholdY = c(4, 6), xlab = "a B-cell", ylab = "a T-cell")
```

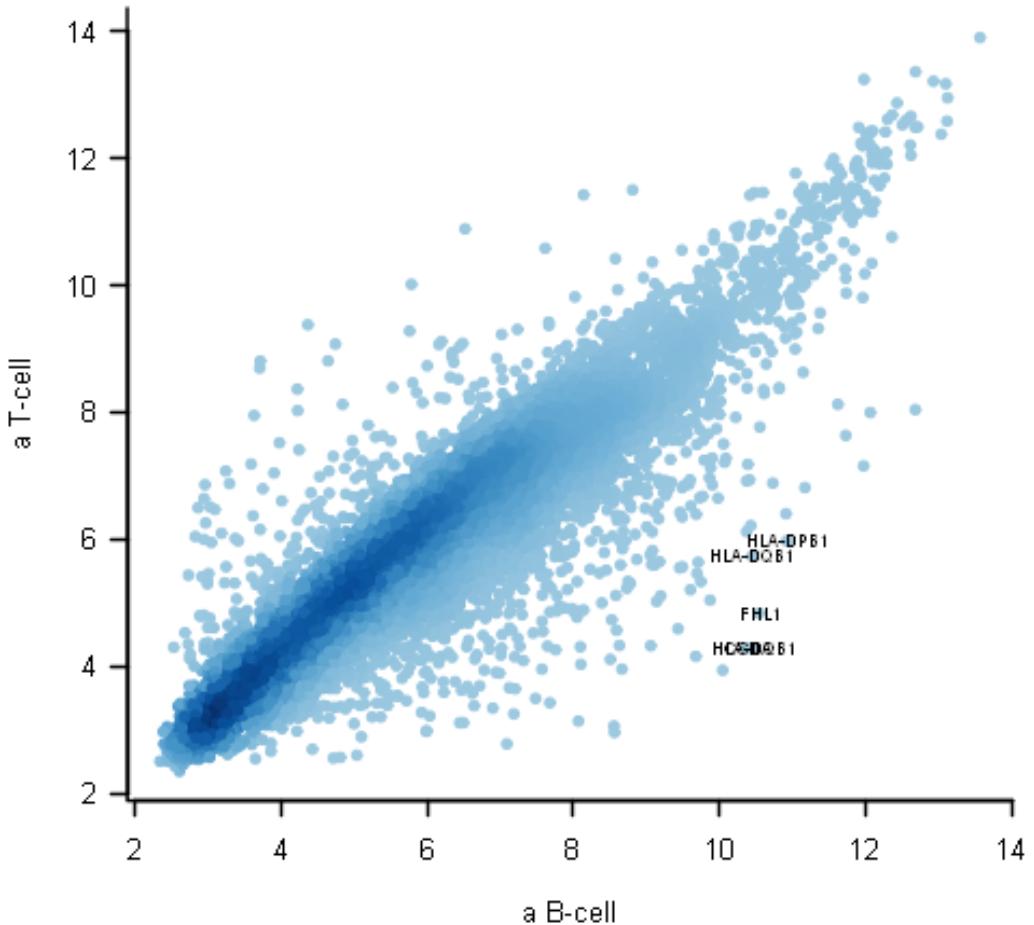


Figure 2: Correlations in gene expression profiles between a B-cell and a T-cell (samples 84004 and 01003). Some potentially interesting genes are indicated by their gene symbol.

One can also show multiple pairwise comparisons in a pairwise scatterplot matrix.

```
R> plotCombMultSamples(exprs(ALL) [, c("84004", "11002", "01003")])
```

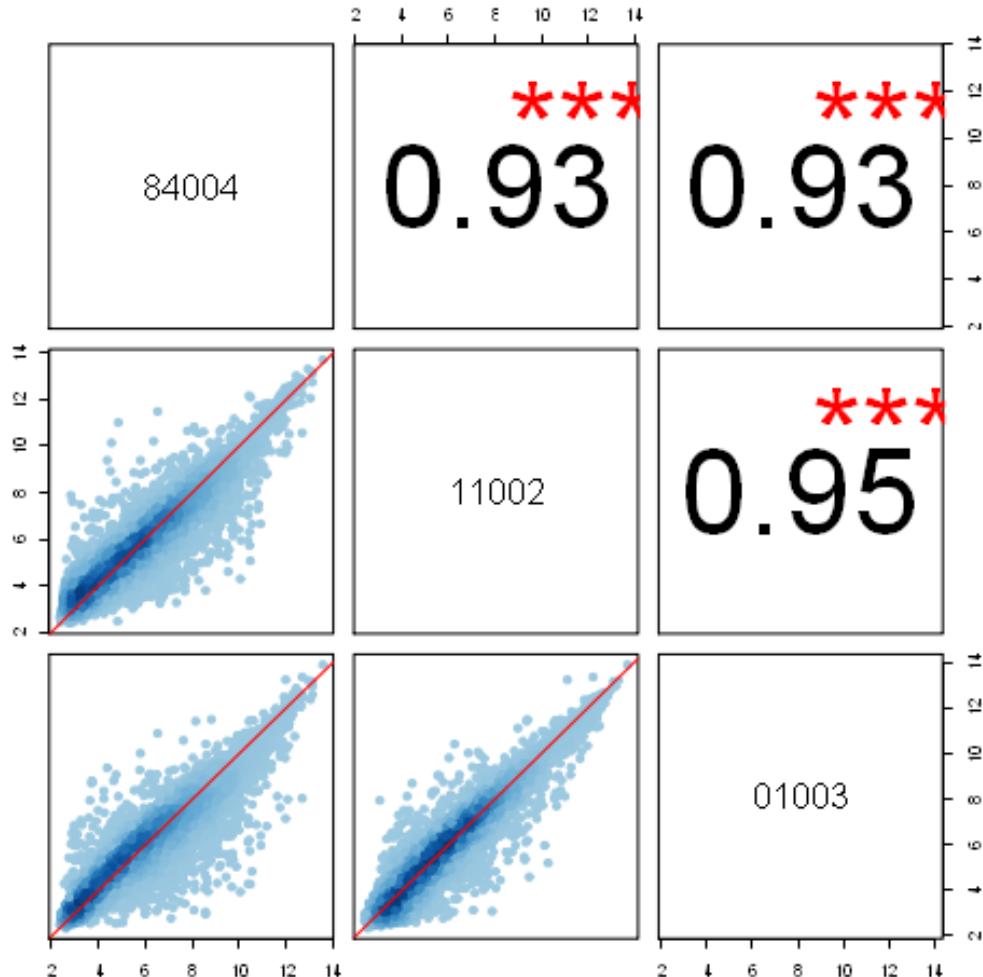


Figure 3: Correlations in gene expression profiles between a B-cell and two T-cell samples (respectively samples 84004, 11002 and 01003).

6.5 Gene lists of log ratios

When analyzing treatments that are primarily interesting relative to a control treatment, it may be of value to look at the log ratios of several treatments (in columns) for a selected list of genes (in rows).

```
R> ALL$BTtype <- as.factor(substr(ALL$BT, 0, 1))
R> ALL2 <- ALL[, ALL$BT != "T1"]
R> ALL2$BTtype <- as.factor(substr(ALL2$BT, 0, 1))
R> tTestResult <- tTest(ALL, "BTtype", probe2gene = FALSE)
R> topGenes <- rownames(tTestResult)[1:20]
R> LogRatioALL <- computeLogRatio(ALL2, reference = list(var = "BT", level = "B"))
R> a <- plotLogRatio(e = LogRatioALL[topGenes, ], openFile = FALSE, tooltipvalues = FALSE, device
```

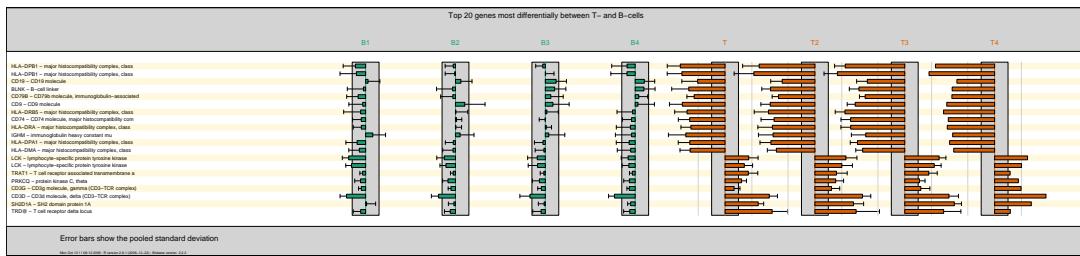


Figure 4: Log ratios of the 20 genes that are most differentially expressed between B-cell and two T-cells.

7 Pathway analysis

7.1 Minus log p

7.2 Gene set enrichment analysis

8 Software used

- R version 2.8.1 (2008-12-22), i386-pc-mingw32
- Locale: LC_COLLATE=English_United States.1252;LC_CTYPE=English_United States.1252;LC_MONETARY=
- Base packages: base, datasets, graphics, grDevices, grid, methods, splines, stats, tools, utils
- Other packages: a4 0.0-18, ALL 1.4.4, annaffy 1.14.0, annotate 1.20.1, AnnotationDbi 1.4.3, Biobase 2.2.2, Cairo 1.4-4, class 7.2-45, cluster 1.11.11, DBI 0.2-4, gdata 2.4.2, genefilter 1.22.0, geneplotter 1.20.0, glmnet 1.1-3, GO.db 2.2.5, gplots 2.6.0, gridSVG 0.5-2, gtools 2.5.0-1, hgu95av2.db 2.2.5, ipred 0.8-7, KEGG.db 2.2.5, lattice 0.17-17, limma 2.16.5, MASS 7.2-45, Matrix 0.999375-23, mlbench 1.1-6, MLInterfaces 1.22.0, mpm 1.0-12, multtest 1.23.3, nlcv 0.1-95, nnet 7.2-45, pamr 1.40.0, plyr 0.1.5, randomForest 4.5-30, RColorBrewer 1.0-2, rda 1.0.1, ROCR 1.0-2, rpart 3.1-42, RSQLite 0.7-1, survival 2.34-1, varSelRF 0.7-1, xtable 1.5-5
- Loaded via a namespace (and not attached): KernSmooth 2.22-22