

# Seamless navigation through combined results of set- & network-based enrichment analysis

**Ludwig Geistlinger**<sup>1</sup>

<sup>1</sup>School of Public Health, City University of New York

**January 5, 2019**

## **Abstract**

The *EnrichmentBrowser* package implements an analysis pipeline for high-throughput gene expression data as measured with microarrays and RNA-seq. In a workflow-like manner, the package brings together a selection of established Bioconductor packages for gene expression data analysis. It integrates a wide range of gene set and network enrichment analysis methods and facilitates combination and exploration of results across methods.

## **Package**

EnrichmentBrowser 2.13.1

Report issues on <https://github.com/lgeistlinger/EnrichmentBrowser/issues>

## Contents

1	Introduction . . . . .	3
2	Reading expression data from file . . . . .	3
3	Types of expression data . . . . .	4
3.1	Microarray data . . . . .	4
3.2	RNA-seq data . . . . .	6
4	Normalization . . . . .	6
5	Differential expression . . . . .	8
6	ID mapping . . . . .	11
7	Enrichment analysis . . . . .	12
7.1	Set-based enrichment analysis . . . . .	12
7.2	Network-based enrichment analysis . . . . .	15
8	Combining results . . . . .	19
9	Putting it all together . . . . .	20
10	Advanced: configuration parameters . . . . .	20
11	For non-R users: command line invocation . . . . .	20
A	A primer on terminology and statistical theory . . . . .	21
A.1	Where does it all come from? . . . . .	21
A.2	Gene sets, pathways & regulatory networks . . . . .	22
A.3	Resources . . . . .	22
A.4	Gene set analysis vs. gene set enrichment analysis . . . . .	22
A.5	Underlying null: competitive vs. self-contained . . . . .	23
A.6	Generations: ora, fcs & topology-based . . . . .	23
B	Frequently asked questions . . . . .	24

# 1 Introduction

---

The *EnrichmentBrowser* package implements essential functionality for the enrichment analysis of gene expression data. The analysis combines the advantages of set-based and network-based enrichment analysis to derive high-confidence gene sets and biological pathways that are differentially regulated in the expression data under investigation. Besides, the package facilitates the visualization and exploration of such sets and pathways.

The following instructions will guide you through an end-to-end expression data analysis workflow including:

1. Preparing the data
2. Preprocessing of the data
3. Differential expression (DE) analysis
4. Defining gene sets of interest
5. Executing individual enrichment methods
6. Combining the results of different methods
7. Visualize and explore the results

All of these steps are modular, i.e. each step can be executed individually and fine-tuned with several parameters. In case you are interested in a particular step, you can directly move on to the respective section. For example, if you have differential expression already calculated for each gene, and you are now interested whether certain gene functions are enriched for differential expression, section *Set-based enrichment analysis* would be the one you should go for. The last section *Putting it all together* also demonstrates how to wrap the whole workflow into a single function, making use of suitably chosen defaults.

# 2 Reading expression data from file

---

Typically, the expression data is not already available in *R* but rather has to be read in from file. This can be done using the function `readSE`, which reads the expression data (`exprs`) along with the phenotype data (`colData`) and feature data (`rowData`) into a *SummarizedExperiment*.

```
library(EnrichmentBrowser)
data.dir <- system.file("extdata", package="EnrichmentBrowser")
exprs.file <- file.path(data.dir, "exprs.tab")
cdat.file <- file.path(data.dir, "colData.tab")
rdat.file <- file.path(data.dir, "rowData.tab")
se <- readSE(exprs.file, cdat.file, rdat.file)
```

The man pages provide details on file format and the *SummarizedExperiment* data structure.

```
?readSE
?SummarizedExperiment
```

## EnrichmentBrowser

*Note:* Previous versions of the *EnrichmentBrowser* used the *ExpressionSet* data structure. The migration to *SummarizedExperiment* in the current release of the *EnrichmentBrowser* is done to reflect recent developments in *Bioconductor*, which discourage use of *ExpressionSet* in favor of *SummarizedExperiment*. Major reasons are the compatibility of *SummarizedExperiment* with operations on genomic regions as well as efficient dealing with big data.

To enable a smooth transition, all functions of the *EnrichmentBrowser* are still accepting also an *ExpressionSet* as input, but are consistently returning a *SummarizedExperiment* as output.

Furthermore, users can always coerce from *SummarizedExperiment* to *ExpressionSet* via

```
eset <- as(se, "ExpressionSet")
```

and vice versa

```
se <- as(eset, "SummarizedExperiment")
```

## 3 Types of expression data

---

The two major data types processed by the *EnrichmentBrowser* are microarray (intensity measurements) and RNA-seq (read counts) data.

Although RNA-seq has become the *de facto* standard for transcriptomic profiling, it is important to know that many methods for differential expression and gene set enrichment analysis have been originally developed for microarray data.

However, differences in data distribution assumptions (microarray: quasi-normal, RNA-seq: negative binomial) made adaptations in differential expression analysis and, to some extent, also in gene set enrichment analysis necessary.

Thus, we consider two example datasets – a microarray and a RNA-seq dataset, and discuss similarities and differences of the respective analysis steps.

### 3.1 Microarray data

To demonstrate the functionality of the package for microarray data, we consider expression measurements of patients with acute lymphoblastic leukemia [1]. A frequent chromosomal defect found among these patients is a translocation, in which parts of chromosome 9 and 22 swap places. This results in the oncogenic fusion gene BCR/ABL created by positioning the ABL1 gene on chromosome 9 to a part of the BCR gene on chromosome 22.

We load the *ALL* dataset

```
library(ALL)
data(ALL)
```

and select B-cell ALL patients with and without the BCR/ABL fusion as described previously [2].

```
ind.bs <- grep("^B", ALL$BT)
ind.mut <- which(ALL$mol.biol %in% c("BCR/ABL", "NEG"))
```

## EnrichmentBrowser

```
sset <- intersect(ind.bs, ind.mut)
all.eset <- ALL[, sset]
```

We can now access the expression values, which are intensity measurements on a log-scale for 12,625 probes (rows) across 79 patients (columns).

```
dim(all.eset)
## Features Samples
## 12625 79

exprs(all.eset)[1:4,1:4]
## 01005 01010 03002 04007
## 1000_at 7.597323 7.479445 7.567593 7.905312
## 1001_at 5.046194 4.932537 4.799294 4.844565
## 1002_f_at 3.900466 4.208155 3.886169 3.416923
## 1003_s_at 5.903856 6.169024 5.860459 5.687997
```

As we often have more than one probe per gene, we summarize gene expression values as the average of the corresponding probe values.

```
allSE <- probe2gene(all.eset)
## Loading required package: hgu95av2.db
## Loading required package: AnnotationDbi
## Loading required package: org.Hs.eg.db
##
##
## Encountered 663 from.IDs with >1 corresponding to.ID
## (the first to.ID was chosen for each of them)
head(rownames(allSE))
## [1] "5595" "7075" "1557" "643" "1843" "4319"
```

Note, that the mapping from probe to gene is done automatically as long as as you have the corresponding annotation package, here the [hgu95av2.db](#) package, installed. Otherwise, the mapping can be defined in the `rowData` slot.

```
rowData(se, use.names=TRUE)
## DataFrame with 767 rows and 2 columns
## PROBEID ENTREZID
## <character> <character>
## 1000_at 1000_at 5595
## 1010_at 1010_at 5600
## 1011_s_at 1011_s_at 7531
## 1013_at 1013_at 4090
## 1018_at 1018_at 7480
## ... ... ...
## 974_at 974_at 9020
## 976_s_at 976_s_at 5594
```

## EnrichmentBrowser

```
## 983_at      983_at      6300
## 993_at      993_at      7297
## 996_at      996_at      2246
```

### 3.2 RNA-seq data

To demonstrate the functionality of the package for RNA-seq data, we consider transcriptome profiles of four primary human airway smooth muscle cell lines in two conditions: control and treatment with dexamethasone [3].

We load the *airway* dataset

```
library(airway)
data(airway)
```

For further analysis, we remove genes with very low read counts and measurements that are not mapped to an ENSEMBL gene ID.

```
airSE <- airway[grepl("^ENSG", rownames(airway)),]
airSE <- airSE[rowMeans(assay(airSE)) > 10,]
dim(airSE)

## [1] 16055      8

assay(airSE)[1:4,1:4]

##           SRR1039508 SRR1039509 SRR1039512 SRR1039513
## ENSG000000000003      679       448       873       408
## ENSG000000000419      467       515       621       365
## ENSG000000000457      260       211       263       164
## ENSG000000000460       60        55        40        35
```

## 4 Normalization

Normalization of high-throughput expression data is essential to make results within and between experiments comparable. Microarray (intensity measurements) and RNA-seq (read counts) data typically show distinct features that need to be normalized for. The function `normalize` wraps commonly used functionality from *limma* for microarray normalization and from *EDASeq* for RNA-seq normalization. For specific needs that deviate from these standard normalizations, the user should always refer to more specific functions/packages.

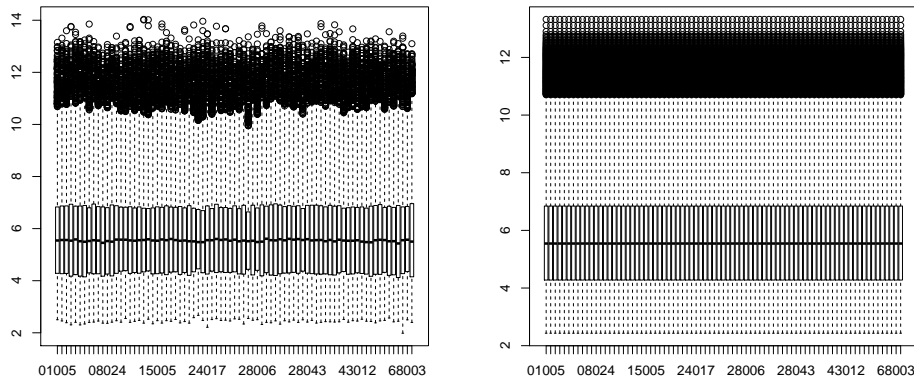
Microarray data is expected to be single-channel. For two-color arrays, it is expected that normalization within arrays has been already carried out, e.g. using `normalizeWithinArrays` from *limma*.

A default quantile normalization based on `normalizeBetweenArrays` from *limma* can be carried out via

```
before.norm <- assay(allSE)
allSE <- normalize(allSE, norm.method="quantile")
after.norm <- assay(allSE)
```

## EnrichmentBrowser

```
par(mfrow=c(1,2))
boxplot(before.norm)
boxplot(after.norm)
```



Note that this is only done for demonstration, as the ALL data has been already RMA-normalized by the authors of the ALL dataset.

RNA-seq data is expected to be raw read counts. Note that normalization for downstream DE analysis, e.g. with [edgeR](#) and [DESeq2](#), is not ultimately necessary (and in some cases even discouraged) as many of these tools implement specific normalization approaches themselves. See the vignette of [EDASeq](#), [edgeR](#), and [DESeq2](#) for details.

In case normalization is desired, between-lane normalization to adjust for sequencing depth can be carried out as demonstrated for microarray data.

```
norm.air <- normalize(airSE, norm.method="quantile")
## Registered S3 method overwritten by 'R.oo':
## method      from
## throw.default R.methodsS3
```

Within-lane normalization to adjust for gene-specific effects such as gene length and GC-content requires to retrieve this information first, e.g. from [BioMart](#) or specific [Bioconductor](#) annotation packages. Both modes are implemented in the [EDASeq](#) function [getGeneLengthAndGCContent](#).

```
ids <- rownames(airSE)
lgc <- EDASeq::getGeneLengthAndGCContent(ids, org="hsa", mode="biomart")
```

Using precomputed information, normalization within and between lanes can be carried out via

```
lgc.file <- file.path(data.dir, "air_lgc.tab")
rowData(airSE) <- read.delim(lgc.file)
norm.air <- normalize(airSE, within=TRUE)
## Normalizing for GC content ...
## Removing 2707 genes due to missing GC content ...
```

## 5 Differential expression

The *EnrichmentBrowser* incorporates established functionality from the *limma* package for differential expression analysis between sample groups. This involves the *voom*-transformation when applied to RNA-seq data. Alternatively, differential expression analysis for RNA-seq data can also be carried out based on the negative binomial distribution with *edgeR* and *DESeq2*.

This can be performed using the function *deAna* and assumes some standardized variable names:

- **GROUP** defines the sample groups being contrasted,
- **BLOCK** defines paired samples or sample blocks, as e.g. for batch effects.

For more information on experimental design, see the *limma user's guide*, chapter 9.

For the ALL dataset, the **GROUP** variable indicates whether the BCR-ABL gene fusion is present (1) or not (0).

```
allSE$GROUP <- ifelse(allSE$mol.biol == "BCR/ABL", 1, 0)
table(allSE$GROUP)

##
##  0  1
## 42 37
```

For the airway dataset, it indicates whether the cell lines have been treated with dexamethasone (1) or not (0).

```
airSE$GROUP <- ifelse(airway$dex == "trt", 1, 0)
table(airSE$GROUP)

##
##  0  1
##  4  4
```

Paired samples, or in general sample batches/blocks, can be defined via a **BLOCK** column in the *colData* slot. For the airway dataset, the sample blocks correspond to the four different cell lines.

```
airSE$BLOCK <- airway$cell
table(airSE$BLOCK)

##
## N052611 N061011 N080611 N61311
##      2      2      2      2
```

For microarray expression data, the *deAna* function carries out a differential expression analysis between the two groups based on functionality from the *limma* package. Resulting fold changes and *t*-test derived *p*-values for each gene are appended to the *rowData* slot.

```
allSE <- deAna(allSE, padj.method="BH")
rowData(allSE, use.names=TRUE)

## DataFrame with 9010 rows and 4 columns
##           FC           limma.STAT           PVAL
##           <numeric>         <numeric>         <numeric>
```



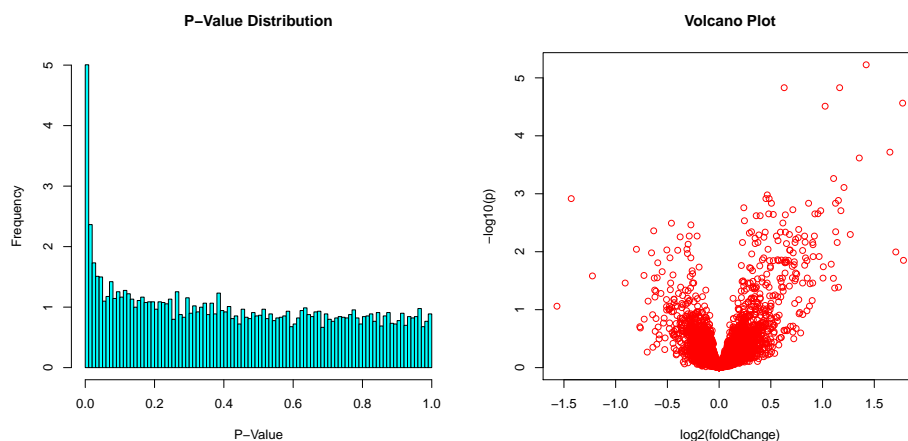
## EnrichmentBrowser

```
## 5595 0.0387427137415798 0.658015072796425 0.512422317574772
## 7075 0.0173335332976947 0.245894327400656 0.806395402373084
## 1557 -0.0507645795671912 -1.27962750198965 0.204384187909288
## 643 -0.0306746085370018 -0.665111958379081 0.507896670249982
## 1843 -0.414988626849719 -1.76920883313719 0.0806797859183739
## ...
## 6300 -0.0447387020187597 -0.93660808882615 0.351786194592616
## 7297 -0.134520145211526 -1.23590025483274 0.220121484483511
## 2246 0.0305107670871747 0.786636249939841 0.433824795604893
## 7850 -0.0209325751219743 -0.239320898612643 0.811470382656448
## 1593 -0.0127671004549253 -0.256145243469366 0.798497864508864
##          ADJ.PVAL
##          <numeric>
## 5595 0.860946231935227
## 7075 0.958384566963236
## 1557 0.683050156652747
## 643  0.860449125663218
## 1843 0.506921109570815
## ...
## 6300 0.782695767830897
## 7297 0.696870897820251
## 2246 0.82259484921151
## 7850 0.958437892107622
## 1593 0.95569417630511
```

Nominal  $p$ -values (PVAL) are corrected for multiple testing (ADJ.PVAL) using the method from Benjamini and Hochberg implemented in the function `p.adjust` from the `stats` package.

To get a first overview, we inspect the  $p$ -value distribution and the volcano plot (fold change against  $p$ -value).

```
par(mfrow=c(1,2))
pdistr(rowData(allSE)$PVAL)
volcano(rowData(allSE)$FC, rowData(allSE)$ADJ.PVAL)
```



The expression change of highest statistical significance is observed for the ENTREZ gene 7525.

## EnrichmentBrowser

```
ind.min <- which.min( rowData(allSE)$ADJ.PVAL )
rowData(allSE, use.names=TRUE)[ ind.min, ]

## DataFrame with 1 row and 4 columns
##           FC          limma.STAT          PVAL
##      <numeric>      <numeric>      <numeric>
## 7525 1.42160480213081 7.01873609978134 6.59964850529937e-10
##           ADJ.PVAL
##      <numeric>
## 7525 5.94628330327473e-06
```

This turns out to be the YES proto-oncogene 1 ([hsa:7525@KEGG](#)).

For RNA-seq data, the `deAna` function carries out a differential expression analysis between the two groups either based on functionality from `limma` (that includes the `voom` transformation), or alternatively, the popular `edgeR` or `DESeq2` package.

Here, we use the analysis based on `edgeR` for demonstration.

```
airSE <- deAna(airSE, de.method="edgeR")
## Excluding 3118 genes not satisfying min.cpm threshold
rowData(airSE, use.names=TRUE)

## DataFrame with 12937 rows and 6 columns
##           length          gc          FC          edgeR.STAT
##      <integer> <numeric>      <numeric>      <numeric>
## ENSG00000000003      8000      0.41 -0.404945626610932      35.8743710016452
## ENSG000000000419    23656      0.398  0.182985434777531      5.90960619951562
## ENSG000000000457    40886      0.403  0.0143477674070905    0.0233923316993606
## ENSG000000000460   190985      0.392 -0.141173372957313    0.492929955080683
## ENSG000000000971    95627      0.352  0.402240426474171    27.8509962017613
## ...
## ENSG00000273270      NA      NA -0.12979385333726    0.901598359265221
## ENSG00000273290      NA      NA  0.505580471641003    23.0905678847793
## ENSG00000273311    2214      0.49  0.00161557580855132  8.04821151395742e-05
## ENSG00000273329      NA      NA -0.222817127090519    1.42723325850574
## ENSG00000273344    2271      0.486  0.0151704005097405    0.005435032737617
##           PVAL          ADJ.PVAL
##      <numeric>      <numeric>
## ENSG00000000003  0.00023480446553515  0.00213458295385677
## ENSG000000000419  0.0388020657296094  0.0915691945173217
## ENSG000000000457  0.88192930278718  0.922279475398735
## ENSG000000000460  0.500971503870371  0.619013213521584
## ENSG000000000971  0.000568781941938381  0.00403820532305421
## ...
## ENSG00000273270  0.367980257149634  0.495892935815196
## ENSG00000273290  0.00106330522558279  0.00639218387702814
## ENSG00000273311  0.993044448477588  0.996356136959404
## ENSG00000273329  0.263765393265588  0.388294594068834
## ENSG00000273344  0.94290685117858  0.962777106053456
```

## 6 ID mapping

Using genomic information from different resources often requires mapping between different types of gene identifiers. Although primary analysis steps such as normalization and differential expression analysis can be carried out independent of the gene ID type, downstream exploration functionality of the *EnrichmentBrowser* is consistently based on NCBI Entrez Gene IDs. It is thus, in this regard, beneficial to initially map gene IDs of a different type to NCBI Entrez IDs.

The function `idTypes` lists the available ID types for the mapping depending on the organism under investigation.

```
idTypes("hsa")
## [1] "ACCNUM"      "ALIAS"      "ENSEMBL"    "ENSEMBLPROT"
## [5] "ENSEMBLTRANS" "ENTREZID"   "ENZYME"     "EVIDENCE"
## [9] "EVIDENCEALL" "GENENAME"   "GO"         "GOALL"
## [13] "IPI"        "MAP"        "OMIM"       "ONTOLOGY"
## [17] "ONTOLOGYALL" "PATH"       "PFAM"       "PMID"
## [21] "PROSITE"    "REFSEQ"     "SYMBOL"     "UCSCKG"
## [25] "UNIGENE"    "UNIPROT"
```

ID mapping for the airway dataset (from ENSEMBL to ENTREZ gene ids) can then be carried out using the function `idMap`.

```
head(rownames(airSE))
## [1] "ENSG000000000003" "ENSG000000000419" "ENSG000000000457" "ENSG000000000460"
## [5] "ENSG000000000971" "ENSG00000001036"

airSE <- idMap(airSE, org="hsa", from="ENSEMBL", to="ENTREZID")
## Encountered 125 from.IDs with >1 corresponding to.ID
## (the first to.ID was chosen for each of them)
## Excluded 1036 from.IDs without a corresponding to.ID
## Encountered 12 to.IDs with >1 from.ID (the first from.ID was chosen for each
of them)

head(rownames(airSE))
## [1] "7105" "8813" "57147" "55732" "3075" "2519"
```

Now, we subject the ALL and the airway gene expression data to the enrichment analysis.

## 7 Enrichment analysis

### 7.1 Set-based enrichment analysis

In the following, we introduce how the *EnrichmentBrowser* package can be used to perform state-of-the-art enrichment analysis of gene sets. We consider the ALL and the airway gene expression data as processed in the previous sections. We are now interested in whether pre-defined sets of genes that are known to work together, e.g. as defined in the Gene Ontology (GO) or the KEGG pathway annotation, are coordinately differentially expressed.

The function `getGenesets` can be used to download gene sets from databases such as GO and KEGG. Here, we use the function to download all KEGG pathways for a chosen organism (here: *Homo sapiens*) as gene sets.

```
kegg.gs <- getGenesets(org="hsa", db="kegg")
```

Analogously, the function `getGenesets` can be used to retrieve GO terms of a selected ontology (here: biological process, BP) as defined in the *GO.db* annotation package.

```
go.gs <- getGenesets(org="hsa", db="go", go.onto="BP", go.mode="GO.db")
```

If provided a file, the function parses user-defined gene sets from GMT file format. Here, we use this functionality for reading a list of already downloaded KEGG gene sets for *Homo sapiens* containing NCBI Entrez Gene IDs.

```
gmt.file <- file.path(data.dir, "hsa_kegg_gs.gmt")
hsa.gs <- getGenesets(gmt.file)
length(hsa.gs)

## [1] 39

hsa.gs[1:2]

## $hsa05416_Viral_myocarditis
## [1] "100509457" "101060835" "1525"      "1604"      "1605"      "1756"
## [7] "1981"      "1982"      "25"        "2534"      "27"        "3105"
## [13] "3106"      "3107"      "3108"      "3109"      "3111"      "3112"
## [19] "3113"      "3115"      "3117"      "3118"      "3119"      "3122"
## [25] "3123"      "3125"      "3126"      "3127"      "3133"      "3134"
## [31] "3135"      "3383"      "3683"      "3689"      "3908"      "4624"
## [37] "4625"      "54205"     "5551"      "5879"      "5880"      "5881"
## [43] "595"       "60"        "637"      "6442"      "6443"      "6444"
## [49] "6445"      "71"        "836"      "841"       "842"      "857"
## [55] "8672"      "940"       "941"      "942"       "958"      "959"
##
## $`hsa04622_RIG-I-like_receptor_signaling_pathway`
## [1] "10010" "1147"  "1432"  "1540"  "1654"  "23586" "26007" "29110"
## [9] "338376" "340061" "3439"  "3440"  "3441"  "3442"  "3443"  "3444"
## [17] "3445"  "3446"  "3447"  "3448"  "3449"  "3451"  "3452"  "3456"
## [25] "3467"  "3551"  "3576"  "3592"  "3593"  "3627"  "3661"  "3665"
## [33] "4214"  "4790"  "4792"  "4793"  "5300"  "54941" "55593" "5599"
## [41] "5600"  "5601"  "5602"  "5603"  "56832" "57506" "5970"  "6300"
## [49] "64135" "64343" "6885"  "7124"  "7186"  "7187"  "7189"  "7706"
```

## EnrichmentBrowser

```
## [57] "79132" "79671" "80143" "841" "843" "8517" "8717" "8737"  
## [65] "8772" "9140" "9474" "9636" "9641" "9755"
```

Currently, the following set-based enrichment analysis methods are supported

```
sbeaMethods()
```

```
## [1] "ora" "safe" "gsea" "gsa" "padog"  
## [6] "globaltest" "roast" "camera" "gsva" "samgs"  
## [11] "ebm" "mgsa"
```

- ORA: Overrepresentation Analysis (simple and frequently used test based on the hypergeometric distribution, see [4] for a critical review),
- SAFE: Significance Analysis of Function and Expression (resampling version of ORA, implements additional test statistics, e.g. Wilcoxon's rank sum, and allows to estimate the significance of gene sets by sample permutation; implemented in the [safe](#) package),
- GSEA: Gene Set Enrichment Analysis (frequently used and widely accepted, uses a Kolmogorov–Smirnov statistic to test whether the ranks of the  $p$ -values of genes in a gene set resemble a uniform distribution [5]),
- PADOG: Pathway Analysis with Down-weighting of Overlapping Genes (incorporates gene weights to favor genes appearing in few pathways versus genes that appear in many pathways; implemented in the [PADOG](#) package),
- ROAST: ROtAtion gene Set Test (uses rotation instead of permutation for assessment of gene set significance; implemented in the [limma](#) and [edgeR](#) packages for microarray and RNA-seq data, respectively),
- CAMERA: Correlation Adjusted MEan RAnk gene set test (accounts for inter-gene correlations as implemented in the [limma](#) and [edgeR](#) packages for microarray and RNA-seq data, respectively),
- GSA: Gene Set Analysis (differs from GSEA by using the maxmean statistic, i.e. the mean of the positive or negative part of gene scores in the gene set; implemented in the [GSA](#) package),
- GSVa: Gene Set Variation Analysis (transforms the data from a gene by sample matrix to a gene set by sample matrix, thereby allowing the evaluation of gene set enrichment for each sample; implemented in the [GSVA](#) package)
- GLOBALTEST: Global testing of groups of genes (general test of groups of genes for association with a response variable; implemented in the [globaltest](#) package),
- SAMGS: Significance Analysis of Microarrays on Gene Sets (extending the SAM method for single genes to gene set analysis [6]),
- EBM: Empirical Brown's Method (combines  $p$ -values of genes in a gene set using Brown's method to combine  $p$ -values from dependent tests; implemented in [Empirical-BrownsMethod](#)),
- MGSA: Model-based Gene Set Analysis (Bayesian modeling approach taking set overlap into account by working on all sets simultaneously, thereby reducing the number of redundant sets; implemented in [mgsa](#)).

See also Appendix A for a comprehensive introduction on underlying statistical concepts.

## EnrichmentBrowser

For demonstration, we perform a basic ORA choosing a significance level  $\alpha$  of 0.1

```
sbea.res <- sbea(method="ora", se=allSE, gs=hsa.gs, perm=0, alpha=0.1)
gsRanking(sbea.res)

## DataFrame with 4 rows and 4 columns
##           GENE.SET NR.GENES NR.SIG.GENES
##           <character> <numeric>      <numeric>
## 1 hsa05130_Pathogenic_Escherichia_coli_infection      43      5
## 2 hsa05206_MicroRNAs_in_cancer      133     10
## 3 hsa04622_RIG-I-like_receptor_signaling_pathway      54      5
## 4 hsa04670_Leukocyte_transendothelial_migration      94      7
##           PVAL
##           <numeric>
## 1      0.0295
## 2      0.0419
## 3      0.0685
## 4      0.0882
```

The result of every enrichment analysis is a ranking of gene sets by the corresponding  $p$ -value. The `gsRanking` function displays only those gene sets satisfying the chosen significance level  $\alpha$ .

While such a ranked list is the standard output of existing enrichment tools, the *Enrichment-Browser* package provides visualization and interactive exploration of resulting gene sets far beyond that point. Using the `eaBrowse` function creates a HTML summary from which each gene set can be inspected in detail (this builds on functionality from the *ReportingTools* package).

The various options are described in Figure 1.

```
eaBrowse(sbea.res)
```

The goal of the *EnrichmentBrowser* package is to provide frequently used enrichment methods. However, it is also possible to exploit its visualization capabilities with user-defined set-based enrichment methods.

This requires to implement a function that takes the characteristic arguments `se` (expression data), `gs` (gene sets), `alpha` (significance level), and `perm` (number of permutations).

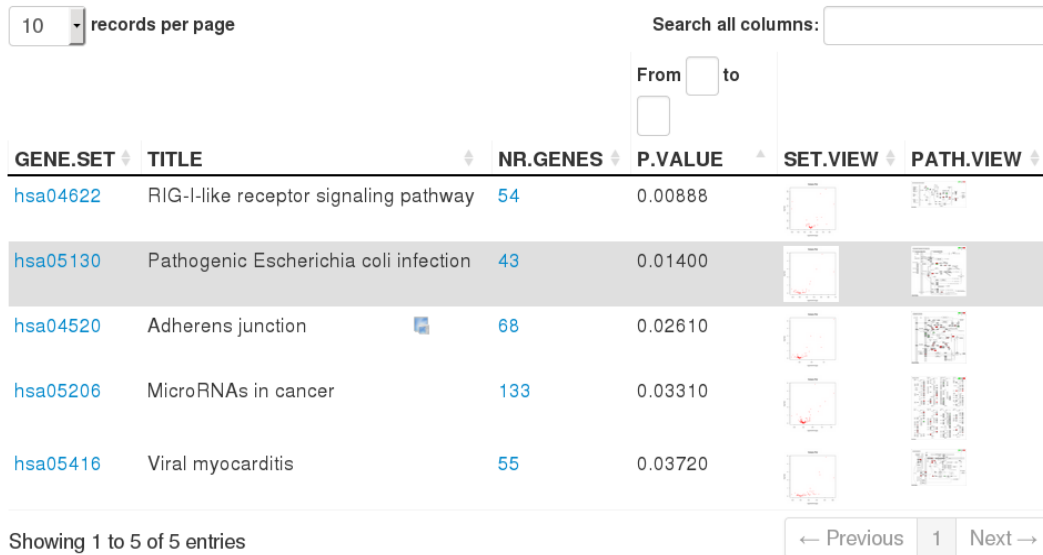
In addition, it must return a numeric vector `ps` storing the resulting  $p$ -value for each gene set in `gs`. The  $p$ -value vector must also be named accordingly (i.e. `names(ps) == names(gs)`).

Let us consider the following dummy enrichment method, which randomly renders five gene sets significant and the remaining insignificant.

```
dummySBEA <- function(se, gs, alpha, perm)
{
  sig.ps <- sample(seq(0,0.05, length=1000),5)
  insig.ps <- sample(seq(0.1,1, length=1000), length(gs)-5)
  ps <- sample(c(sig.ps, insig.ps), length(gs))
  names(ps) <- names(gs)
  return(ps)
}
```

We can plug this method into `sbea` as before.

## ORA - Table of Results



**Figure 1: ORA result view**

For each significant gene set in the ranking, the user can select to view (1) a gene report, that lists all genes of a set along with fold change and DE  $p$ -value, (2) interactive overview plots such as heatmap,  $p$ -value distribution, and volcano plot, (3) the pathway in KEGG with differentially expressed genes highlighted in red.

```
sbea.res2 <- sbea(method=dummySBEA, se=allSE, gs=hsa.gs)
gsRanking(sbea.res2)

## DataFrame with 5 rows and 2 columns
##
##           GENE.SET      PVAL
##           <character> <numeric>
## 1           hsa04520_Adherens_junction  0.00836
## 2 hsa00053_Ascorbate_and_aldarate_metabolism  0.00886
## 3 hsa05410_Hypertrophic_cardiomyopathy_(HCM)  0.0151
## 4           hsa05150_Staphylococcus_aureus_infection  0.0205
## 5 hsa05130_Pathogenic_Escherichia_coli_infection  0.032
```

## 7.2 Network-based enrichment analysis

Having found sets of genes that are differentially regulated in the ALL data, we are now interested whether these findings can be supported by known regulatory interactions.

For example, we want to know whether transcription factors and their target genes are expressed in accordance to the connecting regulations (activation/inhibition). Such information is usually given in a gene regulatory network derived from specific experiments or compiled from the literature ([7] for an example).

## EnrichmentBrowser

There are well-studied processes and organisms for which comprehensive and well-annotated regulatory networks are available, e.g. the RegulonDB for *E. coli* and YeastRACT for *S. cerevisiae*. However, there are also cases where such a network is missing. A basic workaround is to compile a network from regulations in pathway databases such as KEGG.

```
hsa.grn <- compileGRN(org="hsa", db="kegg")
head(hsa.grn)
```

```
##      FROM    TO      TYPE
## [1,] "10000" "100132074" "- "
## [2,] "10000" "1026"      "+"
## [3,] "10000" "1026"      "- "
## [4,] "10000" "1027"      "- "
## [5,] "10000" "10488"     "+"
## [6,] "10000" "107"       "+"
```

Now, we are able to perform enrichment analysis using the compiled network. Currently, the following network-based enrichment analysis methods are supported

```
nbeaMethods()
```

```
## [1] "ggea"      "spia"      "pathnet"   "degraph"   "ganpa"
## [6] "cepa"      "topologygsa" "netgsa"
```

- GGEA: Gene Graph Enrichment Analysis (evaluates consistency of known regulatory interactions with the observed expression data [8]),
- SPIA: Signaling Pathway Impact Analysis (combines ORA with the probability that expression changes are propagated across the pathway topology; implemented in the [SPIA](#) package),
- PathNet: Pathway Analysis using Network Information (applies ORA on combined evidence for the observed signal for gene nodes and the signal implied by connected neighbors in the network; implemented in the [PathNet](#) package),
- DEGraph: Differential expression testing for gene graphs (multivariate testing of differences in mean incorporating underlying graph structure; implemented in the [DEGraph](#) package),
- TopologyGSA: Topology-based Gene Set Analysis (uses Gaussian graphical models to incorporate the dependence structure among genes as implied by pathway topology; implemented in the [topologyGSA](#) package),
- GANPA: Gene Association Network-based Pathway Analysis (incorporates network-derived gene weights in the enrichment analysis; implemented in the [GANPA](#) package),
- CePa: Centrality-based Pathway enrichment (incorporates network centralities as node weights mapped from differentially expressed genes in pathways; implemented in the [CePa](#) package),
- NetGSA: Network-based Gene Set Analysis (incorporates external information about interactions among genes as well as novel interactions learned from data; implemented in the [NetGSA](#) package),

For demonstration, we perform GGEA using the compiled KEGG regulatory network.



## EnrichmentBrowser

```
nbea.res <- nbea(method="ggea", se=allSE, gs=hsa.gs, grn=hsa.grn)
gsRanking(nbea.res)

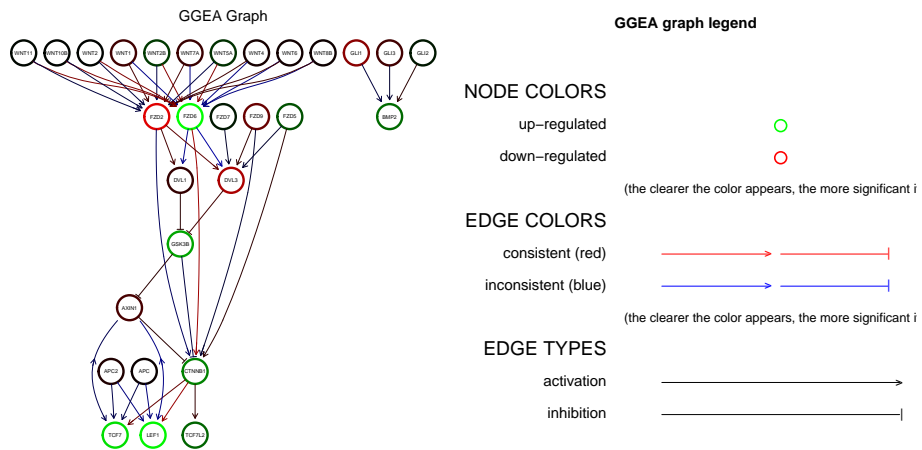
## DataFrame with 7 rows and 5 columns
##           GENE.SET      NR.RELS RAW.SCORE
##           <character> <numeric> <numeric>
## 1          hsa05416_Viral_myocarditis      7      3.29
## 2 hsa04622_RIG-I-like_receptor_signaling_pathway     37     13.9
## 3          hsa04390_Hippo_signaling_pathway     63     22.4
## 4          hsa05217_Basal_cell_carcinoma     17      6.58
## 5          hsa05134_Legionellosis     20      7.38
## 6          hsa04210_Apoptosis     54     18.8
## 7          hsa04520_Adherens_junction     12      4.49
##   NORM.SCORE      PVAL
##   <numeric> <numeric>
## 1      0.47      0.002
## 2      0.377     0.002
## 3      0.356     0.003
## 4      0.387    0.00699
## 5      0.369     0.015
## 6      0.348     0.019
## 7      0.374     0.042
```

The resulting ranking lists, for each statistically significant gene set, the number of relations of the network involving a member of the gene set under study (NR.RELS), the sum of consistencies over the relations of the set (RAW.SCORE), the score normalized by induced network size (NORM.SCORE = RAW.SCORE / NR.RELS), and the statistical significance of each gene set based on a permutation approach.

A GGEA graph for a gene set depicts the consistency of each interaction in the set. Nodes (genes) are colored according to expression (up-/down-regulated) and edges (interactions) are colored according to consistency, i.e. how well the interaction type (activation/inhibition) is reflected in the correlation of the observed expression of both interaction partners.

```
par(mfrow=c(1,2))
ggeaGraph(
  gs=hsa.gs[["hsa05217_Basal_cell_carcinoma"]],
  grn=hsa.grn, se=allSE)
ggeaGraphLegend()
```

## EnrichmentBrowser



As described in the previous section, it is also possible to plug in user-defined network-based enrichment methods.

## 8 Combining results

Different enrichment analysis methods usually result in different gene set rankings for the same dataset. To compare results and detect gene sets that are supported by different methods, the *EnrichmentBrowser* package allows to combine results from the different set-based and network-based enrichment analysis methods. The combination of results yields a new ranking of the gene sets under investigation by specified ranking criteria, e.g. the average rank across methods. We consider the ORA result and the GGEA result from the previous sections and use the function `combResults`.

```
res.list <- list(sbea.res, nbea.res)
comb.res <- combResults(res.list)
```

The combined result can be detailedly inspected as before and interactively ranked as depicted in Figure 2.

```
eaBrowse(comb.res, graph.view=hsa.grn, nr.show=5)
```

### COMB - Table of Results

10 records per page

Search all columns:

From  to  From  to  From  to  From  to  From  to

GENE.SET	TITLE	NR.GENES	ORA.RANK	GGEA.RANK	AVG.RANK	ORA.PVAL	GGEA.PVAL	SET.VIEW	PATH.VIEW	GRAPH.VIEW
hsa05416	Viral myocarditis	55	5	2	3	0.0372	0.006			
hsa05130	Pathogenic Escherichia coli infection	43	2	14	8	0.0140	0.243			
hsa04514	Cell adhesion molecules (CAMs)	107	18	4	11	0.3050	0.010			
hsa04520	Adherens junction	68	3	19	11	0.0261	0.498			
hsa05144	Malaria	45	15	7	11	0.1990	0.067			

Showing 1 to 5 of 5 entries

← Previous 1 Next →

**Figure 2: Combined result view**

By clicking on one of the columns (ORA.RANK, ..., GGEA.PVAL) the result can be interactively ranked according to the selected criterion.

## 9 Putting it all together

There are cases where it is necessary to perform certain steps of the demonstrated enrichment analysis pipeline individually. However, it is often more convenient to run the complete standardized pipeline. This can be done using the all-in-one wrapper function `ebrowser`. For example, the result page displayed in Figure 2 can also be produced from scratch via

```
ebrowser( meth=c("ora", "ggea"),
          exprs=exprs.file, cdat=cdat.file, rdat=rdat.file,
          org="hsa", gs=hsa.gs, grn=hsa.grn, comb=TRUE, nr.show=5)
```

## 10 Advanced: configuration parameters

Similar to *R*'s options settings, the *EnrichmentBrowser* uses certain package-wide configuration parameters, which affect the way in which analysis is carried out and how results are displayed. The settings of these parameters can be examined and, to some extent, also changed using the function `configEBrowser`. For instance, the default directory where the *EnrichmentBrowser* writes results to can be updated via

```
configEBrowser(key="OUTDIR.DEFAULT", value="/my/out/dir")
```

and examined via

```
configEBrowser("OUTDIR.DEFAULT")
## [1] "/my/out/dir"
```

Note that changing these defaults should be done with care, as inappropriate settings might impair the package's functionality. The complete list of incorporated configuration parameters along with their default settings can be inspected via

```
?configEBrowser
```

## 11 For non-R users: command line invocation

The package source contains two scripts in `inst/scripts` to invoke the *EnrichmentBrowser* from the command line using *Rscript*.

The `de_rseq.R` script is a lightweight wrapper script to carry out differential expression analysis of RNA-seq data either based on *limma* (using the `voom`-transformation), *edgeR*, or *DESeq2*.

The `eBrowserCMD.R` implements the full functionality and allows to carry out the various enrichment methods and to produce HTML reports for interactive exploration of results.

The `inst/scripts` folder also contains a README file that comprehensively documents the usage of both scripts.

## A A primer on terminology and statistical theory

### A.1 Where does it all come from?

Test whether known biological functions or processes are over-represented (= enriched) in an experimentally-derived gene list, e.g. a list of differentially expressed (DE) genes. See [4] for a critical review.

Example: Transcriptomic study, in which 12,671 genes have been tested for differential expression between two sample conditions and 529 genes were found DE.

Among the DE genes, 28 are annotated to a specific functional gene set, which contains in total 170 genes. This setup corresponds to a  $2 \times 2$  contingency table,

```
deTable <-
  matrix(c(28, 142, 501, 12000),
        nrow = 2,
        dimnames = list(c("DE", "Not.DE"),
                        c("In.gene.set", "Not.in.gene.set")))
deTable
##           In.gene.set Not.in.gene.set
## DE              28             501
## Not.DE          142            12000
```

where the overlap of 28 genes can be assessed based on the hypergeometric distribution. This corresponds to a one-sided version of Fisher's exact test, yielding here a significant enrichment.

```
fisher.test(deTable, alternative = "greater")
##
## Fisher's Exact Test for Count Data
##
## data:  deTable
## p-value = 4.088e-10
## alternative hypothesis: true odds ratio is greater than 1
## 95 percent confidence interval:
##  3.226736      Inf
## sample estimates:
## odds ratio
##  4.721744
```

This basic principle is at the foundation of major public and commercial enrichment tools such as *DAVID* (<https://david.ncifcrf.gov>) and *Pathway Studio* (<https://www.pathwaystudio.com>).

Although gene set enrichment methods have been primarily developed and applied on transcriptomic data, they have recently been modified, extended and applied also in other fields of genomic and biomedical research. This includes novel approaches for functional enrichment analysis of proteomic and metabolomic data as well as genomic regions and disease phenotypes [9, 10, 11, 12].

### A.2 Gene sets, pathways & regulatory networks

*Gene sets* are simple lists of usually functionally related genes without further specification of relationships between genes.

*Pathways* can be interpreted as specific gene sets, typically representing a group of genes that work together in a biological process. Pathways are commonly divided in metabolic and signaling pathways. Metabolic pathways such as glycolysis represent biochemical substrate conversions by specific enzymes. Signaling pathways such as the MAPK signaling pathway describe signal transduction cascades from receptor proteins to transcription factors, resulting in activation or inhibition of specific target genes.

*Gene regulatory networks* describe the interplay and effects of regulatory factors (such as transcription factors and microRNAs) on the expression of their target genes.

### A.3 Resources

*GO* (<http://www.geneontology.org>) and *KEGG* (<http://www.genome.jp/kegg>) annotations are most frequently used for the enrichment analysis of functional gene sets. Despite an increasing number of gene set and pathway databases, they are typically the first choice due to their long-standing curation and availability for a wide range of species.

The Gene Ontology (GO) consists of three major sub-ontologies that classify gene products according to molecular function (MF), biological process (BP) and cellular component (CC). Each ontology consists of GO terms that define MFs, BPs or CCs to which specific genes are annotated. The terms are organized in a directed acyclic graph, where edges between the terms represent relationships of different types. They relate the terms according to a parent-child scheme, i.e. parent terms denote more general entities, whereas child terms represent more specific entities.

The Kyoto Encyclopedia of Genes and Genomes (KEGG) is a collection of manually drawn pathway maps representing molecular interaction and reaction networks. These pathways cover a wide range of biochemical processes that can be divided in 7 broad categories: metabolism, genetic and environmental information processing, cellular processes, organismal systems, human diseases, and drug development. Metabolism and drug development pathways differ from pathways of the other 5 categories by illustrating reactions between chemical compounds. Pathways of the other 5 categories illustrate molecular interactions between genes and gene products.

### A.4 Gene set analysis vs. gene set enrichment analysis

The two predominantly used enrichment methods are:

- Overrepresentation analysis (ORA), testing whether a gene set contains disproportional many genes of significant expression change, based on the procedure outlined in section A.1,
- Gene set enrichment analysis (GSEA), testing whether genes of a gene set accumulate at the top or bottom of the full gene vector ordered by direction and magnitude of expression change [5].

## EnrichmentBrowser

However, the term *gene set enrichment analysis* nowadays subsumes a general strategy implemented by a wide range of methods [13]. Those methods have in common the same goal, although approach and statistical model can vary substantially [4, 14].

To better distinguish from the specific method, some authors use the term *gene set analysis* to denote the general strategy. However, there is also a specific method of this name [15].

### A.5 Underlying null: competitive vs. self-contained

Goeman and Buehlmann, 2007, classified existing enrichment methods into *competitive* and *self-contained* based on the underlying null hypothesis [4].

- *Competitive* null hypothesis: the genes in the set of interest are at most as often DE as the genes not in the set,
- *Self-contained* null hypothesis: no genes in the set of interest are DE.

Although the authors argue that a self-contained null is closer to the actual question of interest, the vast majority of enrichment methods is competitive.

Goeman and Buehlmann further raise several critical issues concerning the  $2 \times 2$  ORA:

- rather arbitrary classification of genes in DE / not DE,
- based on gene sampling, although sampling of subjects is appropriate,
- unrealistic independence assumption between genes, resulting in highly anti-conservative  $p$ -values.

With regard to these statistical concerns, GSEA is considered superior:

- takes all measured genes into account,
- subject sampling via permutation of class labels,
- the incorporated permutation procedure implicitly accounts for correlations between genes.

However, the simplicity and general applicability of ORA is unmet by subsequent methods improving on these issues. For instance, GSEA requires the expression data as input, which is not available for gene lists derived from other experiment types. On the other hand, the involved sample permutation procedure has been proven inaccurate and time-consuming [15, 16, 17].

### A.6 Generations: ora, fcs & topology-based

Khatri *et al.*, 2012, have taken a slightly different approach by classifying methods along the timeline of development into three generations [14]:

1. Generation: ORA methods based on the  $2 \times 2$  contingency table test,
2. Generation: functional class scoring (FCS) methods such as GSEA, which compute gene set (= functional class) scores by summarizing per-gene DE statistics,
3. Generation: topology-based methods, explicitly taking into account interactions between genes as defined in signaling pathways and gene regulatory networks ([8] for an example).

## EnrichmentBrowser

Although topology-based (also: network-based) methods appear to be most realistic, their straightforward application can be impaired by features that are not detectable on the transcriptional level (such as protein-protein interactions) and insufficient network knowledge [7, 18].

Given the individual benefits and limitations of existing methods, cautious interpretation of results is required to derive valid conclusions. Whereas no single method is best suited for all application scenarios, applying multiple methods can be beneficial. This has been shown to filter out spurious hits of individual methods, thereby reducing the outcome to gene sets accumulating evidence from different methods [19, 20].

## B Frequently asked questions

---

### 1. How to cite the *EnrichmentBrowser*?

Geistlinger L, Csaba G and Zimmer R. Bioconductor's EnrichmentBrowser: seamless navigation through combined results of set- & network-based enrichment analysis. *BMC Bioinformatics*, **17**:45, 2016.

### 2. Is it possible to apply the *EnrichmentBrowser* to simple gene lists?

Enrichment methods implemented in the *EnrichmentBrowser* are, except for ORA, expression-based (and also draw their strength from that). The set-based methods GSEA, SAFE, and SAMGS use sample permutation, involving recomputation of differential expression, for gene set significance estimation, i.e. they require the complete expression matrix. The network-based methods require measures of differential expression such as fold change and  $p$ -value to score interactions of the network. In addition, visualization of enriched gene sets is explicitly designed for expression data. Thus, for simple gene list enrichment, tools like *DAVID* (<https://david.ncifcrf.gov>) and *GeneAnalytics* (<http://geneanalytics.genecards.org>) are more suitable, and it is recommended to use them for this purpose.

## References

- [1] Chiaretti S, Li X, Gentleman R, Vitale A, Vignetti M, et al. Gene expression profile of adult T-cell acute lymphocytic leukemia identifies distinct subsets of patients with different response to therapy and survival. *Blood*, 103(7):2771–8, 2004.
- [2] Gentleman R, Carey V, Huber W, Irizarry R, and Dudoit S. Bioinformatics and computational biology solutions using R and Bioconductor. *Springer*, New York, 2005.
- [3] Himes BE, Jiang X, Wagner P, Hu R, Wang Q, et al. RNA-Seq transcriptome profiling identifies CRISPLD2 as a glucocorticoid responsive gene that modulates cytokine function in airway smooth muscle cells. *PLoS One*, 9(6):e99625, 2014.
- [4] Goeman JJ and Buehlmann P. Analyzing gene expression data in terms of gene sets: methodological issues. *Bioinformatics*, 23(8):980–7, 2007.



## EnrichmentBrowser

- [5] Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci USA*, 102(43):15545–50, 2005.
- [6] Dinu I, Potter JD, Mueller T, Liu Q, Adewale AJ, et al. Improving gene set analysis of microarray data by SAM-GS. *BMC Bioinformatics*, 8:242, 2007.
- [7] Geistlinger L, Csaba G, Dirmeier S, Kueffner R, and Zimmer R. A comprehensive gene regulatory network for the diauxic shift in *Saccharomyces cerevisiae*. *Nucleic Acids Res*, 41(18):8452–63, 2013.
- [8] Geistlinger L, Csaba G, Kueffner R, Mulder N, and Zimmer R. From sets to graphs: towards a realistic enrichment analysis of transcriptomic systems. *Bioinformatics*, 27(13):i366–73, 2011.
- [9] Lavalley-Adam M and Yates JR. Using PSEA-Quant for protein set enrichment analysis of quantitative mass spectrometry-based proteomics. *Curr Protoc Bioinformatics*, 53:13.28.1–16, 2016.
- [10] Chagoyen M, Lopez-Ibanez J, and Pazos F. Functional analysis of metabolomics data. *Methods Mol Biol*, 1415:399–406, 2016.
- [11] McLean CY, Bristor D, Hiller M, Clarke SL, Schaar BT, et al. GREAT improves functional interpretation of cis-regulatory regions. *Nat Biotechnol*, 28(5):495–501, 2010.
- [12] Ried JS, Döring A, Oexle K, Meisinger C, Winkelmann J, et al. PSEA: Phenotype set enrichment analysis—a new method for analysis of multiple phenotypes. *Genet Epidemiol*, 36(3):244–52, 2012.
- [13] Huang da W, Sherman BT, and Lempicki RA. Bioinformatics enrichment tools: paths toward the comprehensive functional analysis of large gene lists. *Nucleic Acids Res*, 37(1):1–13, 2009.
- [14] Khatri P, Sirota M, and Butte AJ. Ten years of pathway analysis: current approaches and outstanding challenges. *PLoS Comput Biol*, 8(2):e1002375, 2012.
- [15] Efron B and Tibshirani R. On testing the significance of sets of genes. *Ann Appl Stat*, 1(1):107–129, 2007.
- [16] Phipson B and Smyth GK. Permutation p-values should never be zero: calculating exact p-values when permutations are randomly drawn. *Stat Appl Genet Mol Biol*, 9:A39, 2010.
- [17] Larson JL and Owen A. Moment based gene set tests. *BMC Bioinformatics*, 16:132, 2015.
- [18] Bayerlova M, Jung K, Kramer F, Klemm F, Bleckmann A, and Beissbarth T. Comparative study on gene set and pathway topology-based enrichment methods. *BMC Bioinformatics*, 16:334, 2015.
- [19] Geistlinger L, Csaba G, and Zimmer R. Bioconductor's EnrichmentBrowser: seamless navigation through combined results of set- & network-based enrichment analysis. *BMC Bioinformatics*, 17:45, 2016.
- [20] Alhamdoosh M, Ng M, Wilson NJ, Sheridan JM, Huynh H, et al. Combining multiple tools outperforms individual methods in gene set enrichment analyses. *Bioinformatics*, 33(3):414–24, 2017.