

# Test association between phenotype and gene expression

Evarist Planet

Bioinformatics & Biostatistics Unit

IRB Barcelona

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Individual gene(s) association with phenotype(s)</b>	<b>2</b>
2.1	Creating an <code>epheno</code> . . . . .	2
2.2	Useful methods for the <code>epheno</code> object . . . . .	4
2.3	Export an <code>epheno</code> . . . . .	7
<b>3</b>	<b>Gene set(s) association with phenotype(s)</b>	<b>7</b>
3.1	Plots that use <code>epheno</code> as input . . . . .	8
3.2	GSEA (Gene Set Enrichment Analysis) . . . . .	10

## 1 Introduction

Imagine a situation where we have gene expression and phenotype variables and we want to test the association of each gene with phenotype. We would probably be interested in testing association of groups of genes (or gene sets) with phenotype. This library provides the tools to do both things in a way that is efficient, structured, fast and scalable. We also provide tools to do *GSEA* (Gene set enrichment analysis) of all phenotype variables at once.

The functions and methods presented on this *vignette* provide tools to easily test association between gene expression levels of individual genes or gene sets of genes and the selected phenotypes of a given gene expression dataset. These can be particularly useful for datasets arising from RNAseq or microarray gene expression studies.

We will load the `ExpressionSet` of a cohort (*GSE2034*) we downloaded from GEO.

```

> options(width=100)
> library(phenoTest)
> data(eset)
> eset

ExpressionSet (storageMode: lockedEnvironment)
assayData: 1000 features, 100 samples
  element names: exprs
protocolData: none
phenoData
  sampleNames: GSM36793 GSM36796 ... GSM36924 (100 total)
  varLabels: PID GEOaccession ... BrainRelapse (7 total)
  varMetadata: labelDescription
featureData: none
experimentData: use 'experimentData(object)'
Annotation: hgu133a

```

For illustration purposes we selected the first 1000 genes and the first 100 samples, created a continuous variable (named *Tumor.size*) and added a new category to the categorical variable *lymph.node.status* to illustrate the functionality of the package.

```

> Tumor.size <- rnorm(ncol(eset),50,2)
> pData(eset) <- cbind(pData(eset),Tumor.size)
> pData(eset)[1:20,'lymph.node.status'] <- 'positive'

```

## 2 Individual gene(s) association with phenotype(s)

### 2.1 Creating an epheno

The **epheno** object will contain the univariate association between a list of phenotype variables and the gene expression from the given **ExpressionSet**. We will use the **ExpressionPhenoTest** function to create the **epheno** object. We will have to tell this function which phenotype variables we want to test and the type of these variables (if they are *ordinal*, *continuous*, *categorical* or *survival* variables). For this purpose we will create a variable called (for instance) **vars2test**. This variable has to be of class **list** with components *continuous*, *categorical*, *ordinal* and *survival* indicating which phenotype variables should be tested. *continuous*, *categorical* and *ordinal* must be character vectors, *survival* a matrix with columns named *time* and *event*. The names must match names in **names(pData(eset))** (being **eset** the **ExpressionSet** of the cohort we are interested in).

```

> head(pData(eset))

```

	PID	GEOaccession	lymph.node.status	Months2Relapse	Relapse	ER.Status	BrainRelapse
	GSM36793	3	GSM36793	positive	101	0	0
	GSM36796	5	GSM36796	positive	118	0	1
	GSM36797	6	GSM36797	positive	9	1	0
	GSM36798	7	GSM36798	positive	106	0	0
	GSM36800	8	GSM36800	positive	37	1	0
	GSM36801	9	GSM36801	positive	125	0	1

	Tumor.size
GSM36793	49.58729
GSM36796	48.76116
GSM36797	47.94115
GSM36798	52.65456
GSM36800	49.85170
GSM36801	47.42516

```
> survival <- matrix(c("Relapse", "Months2Relapse"), ncol=2, byrow=TRUE)
> colnames(survival) <- c('event', 'time')
> vars2test <- list(survival=survival, categorical='lymph.node.status', continuous='Tumor.size')
> vars2test
```

```
$survival
  event    time
[1,] "Relapse" "Months2Relapse"
```

```
$categorical
[1] "lymph.node.status"
```

```
$continuous
[1] "Tumor.size"
```

Now we have everything we need to create the `epheno` object:

```
> epheno <- ExpressionPhenoTest(eset, vars2test, p.adjust.method='none')
```

```
Performing analysis for continuous variable Tumor.size
Performing analysis for categorical variable lymph.node.status
Performing analysis for survival variable Relapse
```

```
> epheno
```

```
Object of class 'epheno'
featureNames: 1007_s_at, 1053_at, 117_at ... (1000) feature(s)
phenoNames: Tumor.size, lymph.node.status, Relapse. (3) phenotype(s)
P-value adjustment method: none
Annotation: hgu133a
Approach: frequentist
```

```
Type "showMethods(classes='epheno')" for a list of ALL methods
```

P values can also be adjusted afterwards:

```
> p.adjust.method(epheno)
```

```
[1] "none"
```

```
> epheno <- pAdjust(epheno, method='BH')
> p.adjust.method(epheno)
```

```
[1] "BH"
```

The **epheno** object extends the **ExpressionSet** object and therefore methods that are available for **ExpressionSets** are also available for **ephenos**.

The effect of both *continuous*, *categorical* and *ordinal* phenotype variables on gene expression levels are tested via **lmFit** from package **limma** (Smyth [2005]). For *ordinal* variables a single coefficient is used to test its effect on gene expression (trend test), which is then used to obtain a P-value. Gene expression effects on *survival* are tested via Cox proportional hazards model (Cox [1972]), as implemented in function **coxph** from package **survival**.

If we want we can compute posterior probabilities instead of pvalues we can set the argument **approach**='bayesian'. The default value is 'frequentist'.

**ExpressionPhenoTest** implements parallel computing via the function **mclapply** from the package **multicore**. Currently **multicore** only operates on Unix systems. If you are a windows user you should set **mc.cores**=1 (the default).

## 2.2 Useful methods for the epheno object

Some of the methods for the **epheno** objects are shown here.

The object can be subseted by phenotype names:

```
> phenoNames(epheno)

[1] "Tumor.size"          "lymph.node.status" "Relapse"

> epheno[, 'Tumor.size']

Object of class 'epheno'
featureNames: 1007_s_at, 1053_at, 117_at ... (1000) feature(s)
phenoNames: Tumor.size. (1) phenotype(s)
P-value adjustment method: BH
Annotation: hgu133a
Approach: frequentist

Type "showMethods(classes='epheno')" for a list of ALL methods

> epheno[,2]

Object of class 'epheno'
featureNames: 1007_s_at, 1053_at, 117_at ... (1000) feature(s)
phenoNames: Tumor.size. (1) phenotype(s)
P-value adjustment method: BH
Annotation: hgu133a
Approach: frequentist

Type "showMethods(classes='epheno')" for a list of ALL methods
```

or by class (class can be ordinal, continuous, categorical or survival):

```
> phenoClass(epheno)

      Tumor.size lymph.node.status      Relapse
      "continuous"      "categorical"      "survival"

> epheno[,phenoClass(epheno)=='survival']

Object of class 'epheno'
featureNames: 1007_s_at, 1053_at, 117_at ... (1000) feature(s)
phenoNames: Relapse. (1) phenotype(s)
P-value adjustment method: BH
Annotation: hgu133a
Approach: frequentist

Type "showMethods(classes='epheno')" for a list of ALL methods
```

`epheno` objects contain information summarizing the association between genes and phenotypes. `getMeans` can be used to obtain the average expression for each group in categorical and ordinal variables, as well as for categorized version of the continuous variables.

```
> head(getMeans(epheno))

      Tumor.size.[45.6,49.3) Tumor.size.[49.3,51.1) Tumor.size.[51.1,54.2]
1007_s_at      11.850961      11.799970      11.739869
1053_at      7.709362      7.623287      7.889880
117_at      7.774214      8.002486      7.923572
121_at      10.247511      10.269498      10.182801
1255_g_at      5.970874      6.032044      5.695821
1294_at      9.365436      9.570391      9.403034

      lymph.node.status.negative lymph.node.status.positive
1007_s_at      11.779693      11.868594
1053_at      7.700153      7.902028
117_at      7.882203      7.965347
121_at      10.217469      10.297187
1255_g_at      5.856978      6.073553
1294_at      9.444534      9.449257
```

Here we see that tumor size has been categorized into 3 groups. The number of categories can be changed with the argument `continuousCategories` in the call to `ExpressionPhenoTest`.

`epheno` objects also contain fold changes and hazard ratios (for survival variables). These can be accessed with `getSummaryDif`, `getFc` and `getHr`.

```
> head(getSummaryDif(epheno))

      Tumor.size.fc.[49.3,51.1) Tumor.size.fc.[51.1,54.2] lymph.node.status.positive.fc
1007_s_at      -1.035976      -1.080045      1.063559
1053_at      -1.061478      1.133291      1.150192
117_at      1.171431      1.109076      1.059324
121_at      1.015357      -1.045875      1.056812
1255_g_at      1.043312      -1.210039      1.161972
1294_at      1.152650      1.026404      1.003280

      Relapse.HR
1007_s_at      -1.299526
1053_at      -1.138198
117_at      1.175962
121_at      1.100616
1255_g_at      1.093210
1294_at      1.006512
```

```
> head(getFc(epheno))
```

	Tumor.size.fc.[49.3,51.1)	Tumor.size.fc.[51.1,54.2]	lymph.node.status.positive.fc
1007_s_at	-1.035976	-1.080045	1.063559
1053_at	-1.061478	1.133291	1.150192
117_at	1.171431	1.109076	1.059324
121_at	1.015357	-1.045875	1.056812
1255_g_at	1.043312	-1.210039	1.161972
1294_at	1.152650	1.026404	1.003280

```
> head(getHr(epheno))
```

	Relapse.HR
1007_s_at	-1.299526
1053_at	-1.138198
117_at	1.175962
121_at	1.100616
1255_g_at	1.093210
1294_at	1.006512

**ExpressionPhenoTest** also computes P-values. **eBayes** from **limma** package is used for continuous, categorical and ordinal phenotypes. A Cox proportional hazards likelihood-ratio test is used for survival phenotypes. P-values can be accessed with **getSignif**. Notice that a single P-value is reported for each phenotype variable. For categorical variables these corresponds to the overall null hypothesis that there are no differences between groups.

```
> head(getSignif(epheno))
```

	Tumor.size	lymph.node.status.positive.pval	Relapse
1007_s_at	0.9391645	0.9022639	0.9801106
1053_at	0.9493761	0.7988976	0.9801106
117_at	0.9258968	0.9225719	0.9801106
121_at	0.9776473	0.8689036	0.9803989
1255_g_at	0.9450779	0.8254344	0.9803989
1294_at	0.9873946	0.9965107	0.9907749

We can also ask for the variables we sent to the **ExpressionPhenoTest** function:

```
> getVars2test(epheno)
```

```
$continuous
[1] "Tumor.size"

$categorical
[1] "lymph.node.status"

$survival
  event      time
[1,] "Relapse" "Months2Relapse"
```

## 2.3 Export an epheno

Functions `export2csv` and `epheno2html` can be used to export to a comma separated value (csv) or an html file. The html file will have useful links to online databases that will provide information about each known gene. For more information about how to use these functions and examples read their help manuals.

## 3 Gene set(s) association with phenotype(s)

Gene sets can be stored in a list object. Each element of the list will contain one gene set. The names of the list will be the names of the gene sets. Here we select genes at random to build our gene sets:

```
> set.seed(777)
> sign1 <- sample(featureNames(eset))[1:20]
> sign2 <- sample(featureNames(eset))[1:50]
> mySignature <- list(sign1,sign2)
> names(mySignature) <- c('My first signature','Another signature')
> mySignature

$`My first signature`
 [1] "201160_s_at" "200964_at"  "200817_x_at" "201465_s_at" "201165_s_at" "1438_at"
 [7] "200815_s_at" "200643_at"  "201414_s_at" "200719_at"   "201198_s_at" "201126_s_at"
[13] "201046_s_at" "201060_x_at" "201327_s_at" "200082_s_at" "200884_at"   "201325_s_at"
[19] "200819_s_at" "200855_at"

$`Another signature`
 [1] "201046_s_at" "201343_at"  "201260_s_at" "200639_s_at" "200991_s_at" "200602_at"
 [7] "201216_at"  "201372_s_at" "201336_at"  "200792_at"   "201241_at"   "200686_s_at"
[13] "200628_s_at" "200013_at"  "200655_s_at" "200660_at"   "200641_s_at" "200674_s_at"
[19] "201023_at"   "200056_s_at" "201158_at"   "200807_s_at" "201320_at"   "200741_s_at"
[25] "200001_at"   "177_at"      "200051_at"   "201174_s_at" "201318_s_at" "200040_at"
[31] "200885_at"   "201200_at"   "201306_s_at" "201433_s_at" "200004_at"   "200980_s_at"
[37] "200613_at"   "201371_s_at" "200929_at"   "201145_at"   "201013_s_at" "200064_at"
[43] "200891_s_at" "201259_s_at" "200657_at"   "201328_at"   "201225_s_at" "201012_at"
[49] "201409_s_at" "201350_at"
```

Gene sets can also be stored in gene set collection objects. From here on all functions have methods for gene sets stored as `lists`, `GeneSets` or `GeneSetCollections`. You can use the one you feel more comfortable with. We will work with `GeneSetCollection`:

```
> library(GSEABase)
> myGeneSetA <- GeneSet(geneIds=sign1, setName='My first signature')
> myGeneSetB <- GeneSet(geneIds=sign2, setName='Another signature')
> mySignature <- GeneSetCollection(myGeneSetA,myGeneSetB)
> mySignature

GeneSetCollection
names: My first signature, Another signature (2 total)
unique identifiers: 201160_s_at, 200964_at, ..., 201350_at (69 total)
types in collection:
  geneIdType: NullIdentifier (1 total)
  collectionType: NullCollection (1 total)
```

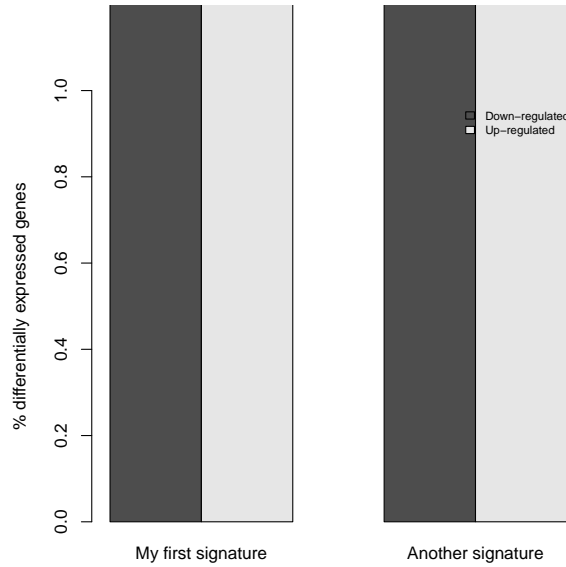


Figure 1: `barplotSignifSignatures`: Number of differentially expressed genes in each gene set that are statistically significant. P-values test for differences in each signature between the number of up and down regulated genes.

### 3.1 Plots that use epheno as input

`barplotSignifSignatures` will plot the percentage of up regulated and down regulated genes that are statistically significant in each signature. In our random selection of genes we did not find any statistically significant genes. Therefore, and just to show the plot we set the alpha value 0.99. The plot can be seen in Figure 1.

```
> barplotSignifSignatures(epheno[, 'lymph.node.status'], mySignature, alpha=0.99, ylim=c(0, 1))
```

By default `barplotSignifSignatures` performs a binomial test (`binom.test` from package `stats`) for each signature to test if the proportions of up regulated and down regulated genes are different. For example, Figure 1 indicates that in the first signature the proportion of up regulated genes is higher than the proportion of down regulated genes. The second signature shows no significant statistical differences.

Sometimes we want to compare the proportions of up and down regulated genes in our signature with the proportions of up and down regulated of all genes in the genome. In this case we may provide a reference signature via



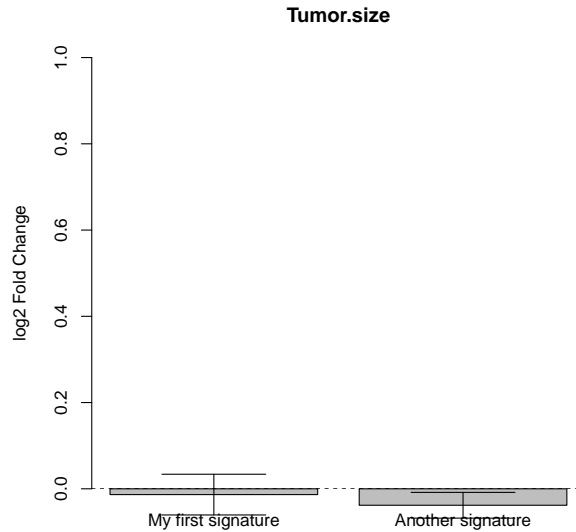


Figure 2: `barplotSignatures`: Average fold change or hazard ratio.

the argument `referenceSignature`. When providing the `referenceSignature` argument a chi-square test comparing the proportion of up and down regulated genes in each signature with the proportion in the reference set will be computed.

When a reference gene set is provided and parameter `testUpDown` is `TRUE` (by default it is `FALSE`) the proportion of up regulated genes is compared with those of the reference gene set. The same is done for down regulated genes.

`barplotSignatures` plots the average log2 fold change or hazard ratio of each phenotype for each gene set. Figure 2 shows an example of it.

```
> barplotSignatures(epheno[, 'Tumor.size'], mySignature, ylim=c(0,1))
```

We can also cluster our samples in two clusters based on the expression levels of one gene set of genes and then test the effect of cluster on phenotypes. For *ordinal* and *continuous* variables a Kruskal-Wallis Rank Sum test is used, for *categorical* variables a chi-square test is used and for *survival* variables a Cox proportional hazards likelihood-ratio test is used. The `heatmapPhenoTest` function can be used to this end. Its results can be seen in Figure 3 and 4.

```
> pvals <- heatmapPhenoTest(eset, mySignature[[1]], vars2test=vars2test[1], heat.kaplan='heat')
```

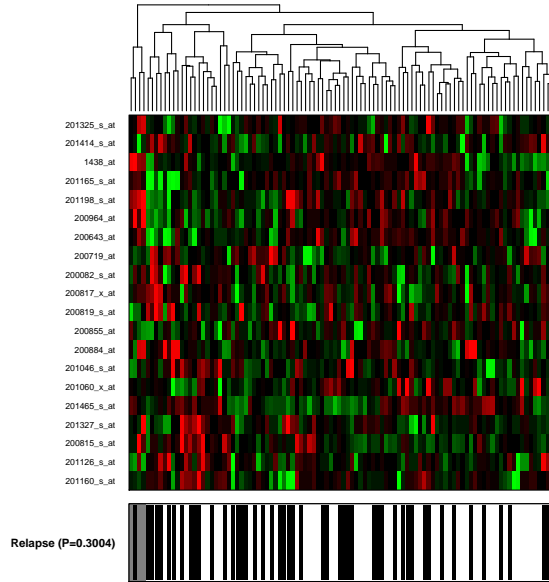


Figure 3: Heatmap produced with `heatmapPhenoTest` function. All variables in `vars2test` that are of class *logical* will be plotted under the heatmap.

```
> pvals

Months2Relapse
"(P=0.3004)"

> pvals <- heatmapPhenoTest(eset,mySignature[[1]],vars2test=vars2test[1],heat.kaplan='kaplan')

> pvals

Months2Relapse
"(P=0.3004)"
```

### 3.2 GSEA (Gene Set Enrichment Analysis)

A popular way to test association between gene sets' gene expression and phenotype is GSEA (Subramanian [2005]). The main idea is to test the association between the gene set *as a whole* and a phenotype.

Although GSEA and several extensions are already available in other *Biconductor* packages, here we implement a slightly different extension. Most GSEA-like approaches assess statistical significance by permuting the values of the phenotype of interest. From a statistical point of view this tests the

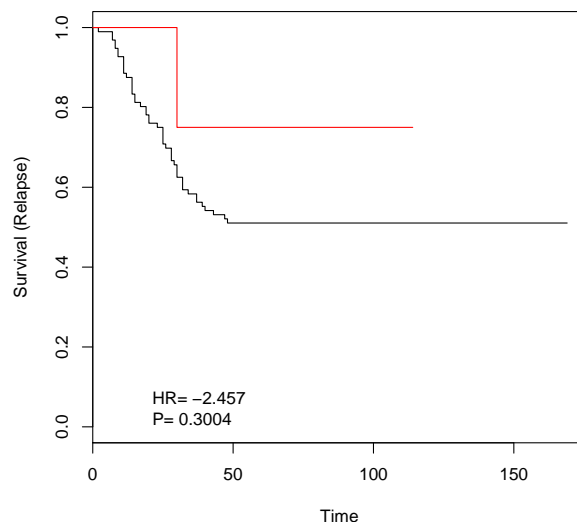


Figure 4: Kaplan-Meier produced with `heatmapPhenoTest` function.

null hypothesis that no genes are associated with phenotype. However in many applications one is actually interested in testing if the proportion of genes associated with phenotype in the gene set is greater than that outside of the gene set. As a simple example, imagine a cancer study where 25% of the genes are differentially expressed. In this setup a randomly chosen gene set will have around 25% of differentially expressed genes, and classical GSEA-like approaches will tend to flag the gene set as statistically significant. In contrast, our implementation will tend to select only gene sets with more than 25% of differentially expressed genes.

We will use the `gsea` method to compute *enrichment scores* (see Subramanian [2005] for details about the enrichment scores) and *simulated enrichment scores* (by permuting the selection of genes). The *simulated enrichment scores* are used to compute P-values and FDR. We can summarize the results obtained using the `summary` method. The following chunk of code is an illustrative example of it:

```
> my.gsea <- gsea(x=epheno,gsets=mySignature,B=1000,p.adjust='BH')

2 gene set(s) were provided and 1000 permutations were assigned,
  therefore 500 permutations will be computed on each gene set.
2 gene set(s) were provided and 1000 permutations were assigned,
  therefore 500 permutations will be computed on each gene set.
2 gene set(s) were provided and 1000 permutations were assigned,
```

```

    therefore 500 permutations will be computed on each gene set.
2 gene set(s) were provided and 1000 permutations were assigned,
    therefore 500 permutations will be computed on each gene set.

> my.gsea

Object of class 'gseaData'
You can use the summary method to produce result summaries.
You can use the getEs, getNes, getEsSim and getFcHr methods to easily access its data.
Gam approximation was not used.
The tested variables are:
  Tumor.size.fc.[49.3,51.1), Tumor.size.fc.[51.1,54.2], lymph.node.status.positive.fc, Relapse.HR
The tested gene sets (for each variable) are:
  My first signature, Another signature

> summary(my.gsea)

      variable      geneSet  n      es      nes  pval.es  pval.nes
1  Tumor.size.fc.[49.3,51.1) My first signature 20 -0.2940722 -0.8809210 0.6212760 0.6484648
2  Tumor.size.fc.[49.3,51.1) Another signature 50 -0.3479237 -1.2787640 0.2241363 0.2927300
3  Tumor.size.fc.[51.1,54.2] My first signature 20  0.2313047  0.6532771 0.9066354 0.9357443
4  Tumor.size.fc.[51.1,54.2] Another signature 50 -0.3117507 -1.1155027 0.5793075 0.6159795
5 lymph.node.status.positive.fc My first signature 20 -0.3768407 -1.1312888 0.2939268 0.2759162
6 lymph.node.status.positive.fc Another signature 50  0.3523689  1.2699886 0.2561693 0.2759162
7 Relapse.HR My first signature 20  0.3184688  1.0656066 0.3706713 0.3603246
8 Relapse.HR Another signature 50 -0.2979001 -1.1443660 0.3706713 0.3603246

      fdr
1 0.6484648
2 0.2927300
3 0.9357443
4 0.3079897
5 0.2759162
6 0.1549552
7 0.3603246
8 0.2559304

```

We receive one message for each phenotype we are testing.

We can produce plots as follows:

```
> plot(my.gsea)
```

This will produce two plots (one for *enrichment score* and another for *normalised enrichment score*) for every phenotype and gene set (in our case 12 plots). Following code shows an example on plotting only *enrichment score* for variable *Relapse* on the first gene set of genes. Plot can be seen in Figure 5.

```
> my.gsea <- gsea(x=epheno[, 'Relapse'], gsets=mySignature[1], B=100, p.adjust='BH')
```

```

1 gene set(s) were provided and 100 permutations were assigned,
    therefore 100 permutations will be computed on each gene set.

```

```
> summary(my.gsea)
```

```

      variable      geneSet  n      es      nes  pval.es  pval.nes      fdr
1 Relapse.HR My first signature 20 0.3184688 1.081041 0.3686867 0.3686856 0.3686856

```

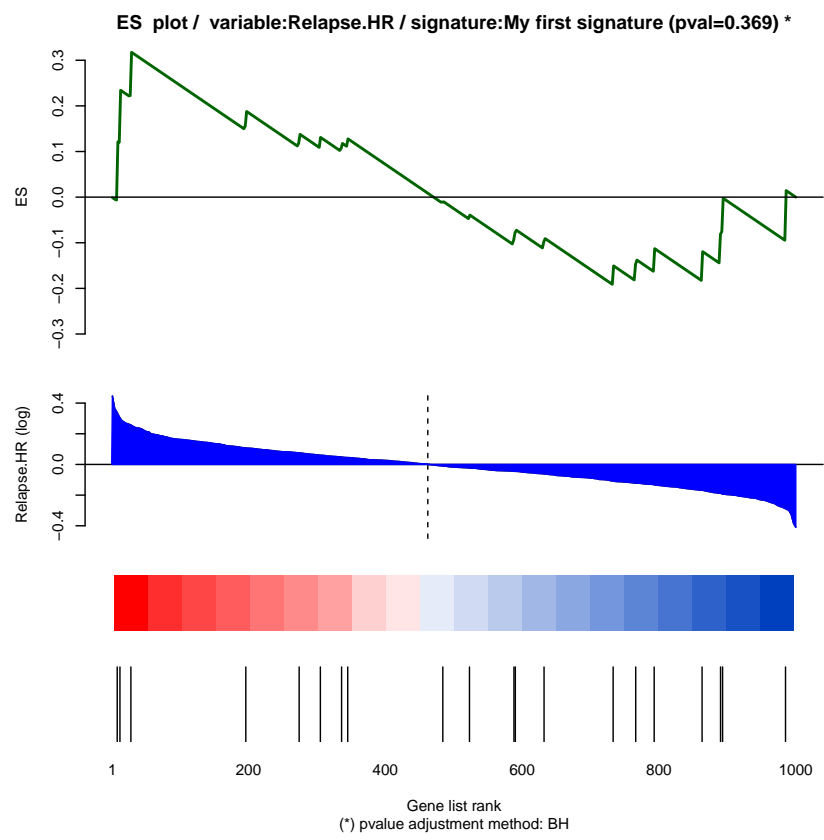


Figure 5: GSEA plot.

```
> plot(my.gsea,es.nes='es',selGsets='My first signature')
```

**gsea** can be used not only with **epheno** objects but also with objects of class **numeric** or **matrix**. For more information read the **gsea** function help.

Following similar ideas to Virtaneva [2001] we also implemented a Wilcoxon test. This can be used instead of the permutation test which can be slow if we use a lot of permutations and we can not use the **multicore** package. The plot we will obtain will also be different. Instead of plotting the *enrichment scores* we will plot the density function and the mean log2 fold change or hazard ratio of the genes that belong to our gene set. This will allow us to compare how similar/different from 0 the mean of our gene set is. The plot using Wilcoxon test can be seen in Figure 6.

```
> my.gsea <- gsea(x=epheno[, 'Relapse'],gsets=mySignature,B=100,test='wilcox',p.adjust='BH')
```

```
2 gene set(s) were provided and 100 permutations were assigned,
  therefore 50 permutations will be computed on each gene set.
```

```
> summary(my.gsea)
```

	variable	geneSet	n	es	pval
1	Relapse.HR	My first signature	20	0.002077225	0.8694878
2	Relapse.HR	Another signature	50	-0.014101939	0.8694878

```
> plot(my.gsea,selGsets='My first signature')
```

Notice that using a Wilcoxon test is conceptually very similar to the average gene set fold change presented in figure 2.

A current limitation of **gseaSignatures** is that it does not consider the existence of dependence between genes in the gene set. This will be addressed in future versions. Nevertheless we believe **gseaSignatures** is useful in that it targets the correct null hypothesis that gene set is as enriched as a randomly selected gene set, opposed to testing that there are no enriched genes in the set as is done in GSEA.

## References

- D.R. Cox. Regression models and life tables. *Journal of the Royal Statistical Society Series B*, 34:187–220, 1972.
- G.K. Smyth. Limma: linear models for microarray data. In R. Gentleman, V. Carey, S. Dudoit, R. Irizarry, and W. Huber, editors, *Bioinformatics and Computational Biology Solutions using R and Bioconductor*, pages 397–420. Springer, New York, 2005.



Figure 6: GSEA plot using Wilcoxon test.

- Aravind Subramanian. Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *PNAS*, 102, 2005.
- K. Virtaneva. Expression profiling reveals fundamental biological differences in acute myeloid leukemia with isolated trisomy 8 and normal cytogenetics. *Proc Natl Acad Sci U S A*, 98(98):1124–1129, January 2001. doi: <http://dx.doi.org/10.1073/pnas.98.3.1124>. URL <http://dx.doi.org/10.1073/pnas.98.3.1124>.