

Differential Expression Meta-Analysis with DExMA package

Juan Antonio Villatoro-García^{1,2} and Pedro Carmona-Sáez*

¹Department of Statistics and Operational Research. University of Granada

²Bioinformatics Unit. GENYO, Centre for Genomics and Oncological Research,

*pedro.carmona@genyo.es

April 26, 2022

Abstract

DExMA (Differential Expression Meta-Analysis) performs all the necessary steps of differential expression meta-analysis considering the possible existence of missing genes. It controls the existence of unmeasured genes by two approaches: considering the genes that are contained in a proportion of datasets (set by the users) or imputing the expression values of the unmeasured genes using the *Knn method* in the space of samples. In addition, it allows to apply quality controls, download GEO datasets and show a graphical representation of the results.

packageVersionDExMA 1.4.0

Contents

| | | |
|-------|--|----|
| 1 | Introduction | 3 |
| 2 | Previous steps: Meta-analysis object | 3 |
| 2.1 | Meta-analysis object creation (objectMA) | 4 |
| 2.2 | Adding a new dataset to the meta-analysis object | 7 |
| 3 | Performing Meta-analysis | 9 |
| 3.1 | Gene annotation and quality controls | 9 |
| 3.1.1 | Setting all the datasets in the same annotation | 9 |
| 3.1.2 | Logarithm transformation | 11 |
| 3.1.3 | Heterogeneity study | 12 |
| 3.2 | Imputing missing genes expression values (optional). | 13 |
| 3.3 | Performing meta-analysis: <i>metaAnalysisDE()</i> | 14 |
| 3.3.1 | Effects size combination results | 15 |
| 3.3.2 | P-value combination results | 16 |
| 3.4 | Visualization of the results: heatmap | 17 |
| 4 | Additional information | 18 |

DExMA package

| | | |
|-----|---|----|
| 4.1 | GEO microarray data download | 18 |
| 4.2 | Using RNA-Seq data | 19 |
| 4.3 | Removing Batch Effects | 19 |
| 4.4 | Calculating Effects size | 21 |
| 4.5 | Calculating Individual P-values | 22 |
| 5 | Session info | 22 |

1 Introduction

DExMA is a package designed to perform gene expression meta-analysis. Gene expression meta-analysis comprises a set of methods that combine the results of several differential expression studies into a single common result [1]. Furthermore, this package has the advantage that it allows to take into account those genes that are not measured in a study. The first approach to control these missing genes, **DExMA** offers two options: apply the k-nearest neighbors method (knn method) in the space of samples to impute the expression values of these missing genes [2]. The second option is to take into account those genes that are contained in at least a certain proportion (set by the user) of datasets. The use of only the genes common to all the datasets could lead to the loss of some genes that are not measured in a single study, which would lead to the loss of information. Due to this fact, **DExMA** package is very useful to work with Microarray data, because this type of data is what usually produces the existence of uncommon genes between datasets with different annotations. However, previously normalized RNA-Seq data can also be used, as well as, both Microarray data and normalized RNA-Seq data at the same time.

DExMA package has implemented methods from the three main types of gene expression meta-analysis [1] meta-analysis based on effects sizes combination and meta-analysis based on p-values combination. Once one of the methods has been applied, **DExMA** package provides a specific and adapted results table. Moreover, this package contains some functions that allows the user to carry out a previous quality control in order to results obtained are more reliable. Finally, this package provides some additional function that, for example, help the user to download public Microarray data from NCBI GEO public database [3] or to visualize the significant genes in a heatmap.

This document gives a tutorial-style introduction to all the steps that must be carry out in order to properly perform gene expression meta-analysis by making use of **DExMA** package.

2 Previous steps: Meta-analysis object

DExMA uses a specific object as input, which is a list of nested lists where each nested list corresponds to a study. This object can be created directly by the users or they can use *createObjectMA()* function to create it.

For the examples that are going to be shown, synthetic data will be used. We load the sample data into our R session.

```
> library(DExMA)
> data("DExMAExampleData")
```

- listMatrixEX: a list of four expression arrays
- listPhenodatas: a list of the four phenodata corresponding to the four expression arrays
- listExpression: a list of four ExpressionSets object. It contains the same information as listMatrixEX and listPhenodatas
- ExpressionSetStudy5: an ExpressionSet object similar to the ExpressionSets objects of listExpression.
- maObjectDif: the meta-analysis object created from the listMatrixEX and listPhenodatas objects.

- `maObject`: the meta-analysis object after setting all the studies in Official Gene Symbol annotation

2.1 Meta-analysis object creation (`objectMA`)

As previously stated, the meta-analysis input in DExMA is a list of nested lists. Each nested list contains two elements:

- A gene expression matrix with genes in rows and samples in columns
- A vector of 0 and 1 indicating the group of each sample. 0 represents reference group (usually controls) and 1 represents experimental group (usually cases).

This object can be created directly by the user or we can make use of `createObjectMA()` function, which creates the `*objectMA*` after indicating how the reference and experimental groups are identified.

`createObjectMA()` function allows to create the object needed to perform meta-analysis. In this case, it is necessary to indicate as input of the function the variables that contain the experimental and reference groups:

- `listEX`: a list of dataframes or matrix (genes in rows and samples in columns). A list of `ExpressionSets` can be used too:

```
> #List of expression matrices
> data("DExMAExampleData")
> ls(listMatrixEX)

[1] "Study1" "Study2" "Study3" "Study4"

> head(listMatrixEX$Study1)

      Sample1 Sample2 Sample3 Sample4
100859927 5.439524 6.253319 2.926444 4.4304023
8086      5.769823 5.971453 1.831349 4.0466288
8212      7.558708 5.957130 2.365252 3.4352889
65985     6.070508 7.368602 2.971158 3.7151784
729522    6.129288 5.774229 3.670696 3.9171749
13        7.715065 7.516471 1.349453 0.3390772

> #List of ExpressionSets
> ls(listExpressionSets)

[1] "Study1" "Study2" "Study3" "Study4"

> listExpressionSets$Study1

ExpressionSet (storageMode: lockedEnvironment)
assayData: 200 features, 4 samples
  element names: exprs
protocolData: none
phenoData
  rowNames: Sample1 Sample2 Sample3 Sample4
  varLabels: condition gender organism race
  varMetadata: labelDescription
featureData: none
```

DExMA package

```
experimentData: use 'experimentData(object)'  
Annotation:
```

- `listPheno`: a list of phenodatas (samples in rows and covariables in columns). If the object `listEX` is a list of `ExpressionSets` this element can be null.

```
> data("DExMAExampleData")  
> #Example of a phenodata object  
> ls(listPhenodatas)  
  
[1] "Study1" "Study2" "Study3" "Study4"  
  
> listPhenodatas$Study1  
  
      condition gender    organism race  
Sample1  Diseased Female Homo Sapiens  AA  
Sample2  Diseased Female Homo Sapiens  AA  
Sample3   Healthy Female Homo Sapiens  AA  
Sample4   Healthy Female Homo Sapiens   H
```

- `namePheno`: a list or vector of the different column names or column positions from the pheno used for performing the comparison between groups. Each element of `namePheno` correspond to its equivalent element in the `listPheno`. (default a vector of 1, all the first columns of each elements of `listPheno` are selected)
- `expGroups`: a list or vector of the group names or positions from `namePheno` variable used as experimental group (cases) to perform the comparison (default a vector of 1, all the first groups are selected).
- `refGroups`: a list or vector of the group names or positions from `namePheno` variable used as reference group (controls) to perform the comparison (default a vector of 2, all the second groups are selected).

It is important to note that if any element does not belong to the experimental or the reference group, that sample is not taken into account in the creation of meta-analysis object.

Here, we have included an example to show how exactly the function is used:

Since this function can be a bit complicated if there are many datasets, we recommend creating a vector to keep the column names of the phenodatas that contains the variable that identifies the groups to compare (*namePheno* argument). Moreover, we should create two others lists to indicate how to identify experimental (cases) and reference (controls) groups in these variables (*expGroups* and *refGroups* arguments).

If we look at the example phenodatas list we have the following four objects:

```
> listPhenodatas$Study1  
  
      condition gender    organism race  
Sample1  Diseased Female Homo Sapiens  AA  
Sample2  Diseased Female Homo Sapiens  AA  
Sample3   Healthy Female Homo Sapiens  AA  
Sample4   Healthy Female Homo Sapiens   H
```

In the "Study1" phenoData, the groups variable is "condition". Experimental group is named as "Diseased" and reference group as "Healthy".

DExMA package

```
> listPhenodatas$Study2
```

| | condition | gender | organism | race |
|----------|-----------|--------|--------------|------|
| Sample5 | Diseased | Female | Homo Sapiens | AA |
| Sample6 | Diseased | Female | Homo Sapiens | AA |
| Sample7 | ill | Male | Homo Sapiens | C |
| Sample8 | Healthy | Female | Homo Sapiens | H |
| Sample9 | control | Female | Homo Sapiens | H |
| Sample10 | control | Male | Homo Sapiens | H |

In the "Study2" phenoData, the groups variable is "condition". Experimental group is named as "Diseased" or "ill" and reference group as "Healthy" or "control"

```
> listPhenodatas$Study3
```

| | state | gender | organism | race |
|----------|----------|--------|--------------|------|
| Sample11 | Diseased | Female | Homo Sapiens | AA |
| Sample12 | Diseased | Female | Homo Sapiens | AA |
| Sample13 | Healthy | Female | Homo Sapiens | AA |
| Sample14 | Healthy | Female | Homo Sapiens | H |

In the "Study3" phenoData, the groups variable is "state". Experimental group is named as "Diseased" and reference group as "Healthy".

```
> listPhenodatas$Study4
```

| | state | gender | organism | race |
|----------|---------|--------|--------------|------|
| Sample15 | ill | Female | Homo Sapiens | AA |
| Sample16 | ill | Female | Homo Sapiens | AA |
| Sample17 | ill | Male | Homo Sapiens | AA |
| Sample18 | control | Female | Homo Sapiens | C |
| Sample19 | control | Female | Homo Sapiens | H |
| Sample20 | control | Male | Homo Sapiens | H |

In this phenoData, the groups variable is "state". Experimental group is named as "ill" and reference group as "control".

We all this information we can create the vector for *namePheno* argument and the two list for *expGroups* and *refGroups*:

```
> phenoGroups = c("condition", "condition", "state", "state")
> phenoCases = list(Study1 = "Diseased", Study2 = c("Diseased", "ill"),
+                   Study3 = "Diseased", Study4 = "ill")
> phenoControls = list(Study1 = "Healthy", Study2 = c("Healthy", "control"),
+                      Study3 = "Healthy", Study4 = "control")
```

Then, we can apply more easily *createObjectMA()* function:

```
> newObjectMA <- createObjectMA(listEX=listMatrixEX,
+                               listPheno = listPhenodatas,
+                               namePheno=phenoGroups,
+                               expGroups=phenoCases,
+                               refGroups = phenoControls)
> #Study 1
```

```

> head(newObjectMA[[1]][[1]])

      Sample1 Sample2 Sample3 Sample4
100859927 5.439524 6.253319 2.926444 4.4304023
8086      5.769823 5.971453 1.831349 4.0466288
8212      7.558708 5.957130 2.365252 3.4352889
65985     6.070508 7.368602 2.971158 3.7151784
729522    6.129288 5.774229 3.670696 3.9171749
13        7.715065 7.516471 1.349453 0.3390772

> newObjectMA[[1]][[2]]

[1] 1 1 0 0

> #Study 2
> head(newObjectMA[[2]][[1]])

      Sample5 Sample6 Sample7 Sample8 Sample9 Sample10
100859927 4.367690 6.648252 5.786897 2.820073 2.688687 3.331881
8086      5.937004 5.914434 7.124650 3.234439 3.228559 4.521235
8212      5.294553 4.092231 5.194514 3.368330 2.750491 5.369983
65985     5.685822 5.900363 5.539580 2.153294 3.166798 3.162232
729522    5.733054 5.553966 6.592146 3.275042 2.659061 2.452524
13        6.153159 5.435909 6.991360 4.860632 3.751608 3.121853

> newObjectMA[[2]][[2]]

[1] 1 1 1 0 0 0

```

The result obtained is the proper object to perform meta-analysis (objectMA).

2.2 Adding a new dataset to the meta-analysis object

It may happen that once the meta-analysis object is created we want to add a new dataset before doing the meta-analysis. **DExMA** provides the *elementObjectMA()* function, which allows the creation of an element of the meta-analysis object. This function contains the following arguments:

- **expressionMatrix**: a dataframe or matrix that containing genes in rows and samples in columns. An *ExpressionSet* object can be used too.
- **pheno**: a data frame or a matrix containing samples in rows and covariates in columns. If NULL (default), pheno is extracted from the *ExpressionSet* object
- **groupPheno**: the column name or position from pheno where experimental group (cases) and reference group (control) are identified
- **expGroups**: a vector of the names or positions from groupPheno variable used as experimental group (cases). By default the first group (character) is taken
- **refGroups**: a vector of the names or positions from groupPheno variable used as reference group (control). By default the second group (character) is taken.

As with the *createObjectMA()* function, if any element does not belong to the experimental or the reference group, that sample is not included in the creation of the object.

Here we provided an example of the use of this function, in which we create an element of the meta-analysis object from the information of the "Study 2" of the *listExpressionSets* object:

DExMA package

```
> data("DExMAExampleData")
> ExpressionSetStudy5

ExpressionSet (storageMode: lockedEnvironment)
assayData: 200 features, 6 samples
  element names: exprs
protocolData: none
phenoData
  rowNames: newSample1 newSample2 ... newSample6 (6 total)
  varLabels: condition gender organism race
  varMetadata: labelDescription
featureData: none
experimentData: use 'experimentData(object)'
Annotation:

> library(Biobase)
> pData(ExpressionSetStudy5)

      condition gender  organism race
newSample1  Diseased Female Homo Sapiens  AA
newSample2  Diseased Female Homo Sapiens  AA
newSample3    ill    Male Homo Sapiens    C
newSample4  Healthy Female Homo Sapiens    H
newSample5  control Female Homo Sapiens    H
newSample6  control  Male Homo Sapiens    H
```

We had to load [Biobase](#) package in order to the information in the *ExpressionSet*. In the *phenoData* we can observe that groups variable is "condition". In addition, Experimental group is name as "Diseased" or "ill" and reference group as "Healthy" or "control"

```
> newElem <-elementObjectMA(expressionMatrix = ExpressionSetStudy5,
+                             groupPheno = "condition",
+                             expGroup = c("Diseased", "ill"),
+                             refGroup = c("Healthy", "control"))
> head(newElem[[1]])

      newSample1 newSample2 newSample3 newSample4 newSample5 newSample6
100859927  5.867690  8.148252  7.286897  4.320073  4.188687  4.831881
8086      7.437004  7.414434  8.624650  4.734439  4.728559  6.021235
8212      6.794553  5.592231  6.694514  4.868330  4.250491  6.869983
65985     7.185822  7.400363  7.039580  3.653294  4.666798  4.662232
729522    7.233054  7.053966  8.092146  4.775042  4.159061  3.952524
13        7.653159  6.935909  8.491360  6.360632  5.251608  4.621853

> head(newElem[[2]])

[1] 1 1 1 0 0 0
```

As we can see, we obtain a list that has the same structure as the elements of the meta-analysis object (*objectMA*). This new element can be added to a *objectMA* that has been created previously:

```
> newObjectMA2 <- newObjectMA
> newObjectMA2[[5]] <- newElem
```



```
> head(newObjectMA2[[5]][[1]])
      newSample1 newSample2 newSample3 newSample4 newSample5 newSample6
100859927  5.867690  8.148252  7.286897  4.320073  4.188687  4.831881
8086      7.437004  7.414434  8.624650  4.734439  4.728559  6.021235
8212      6.794553  5.592231  6.694514  4.868330  4.250491  6.869983
65985     7.185822  7.400363  7.039580  3.653294  4.666798  4.662232
729522    7.233054  7.053966  8.092146  4.775042  4.159061  3.952524
13        7.653159  6.935909  8.491360  6.360632  5.251608  4.621853

> newObjectMA2[[5]][[2]]
[1] 1 1 1 0 0 0
```

Moreover, an advantage of this function is that it can be used to create one by one all the elements and finally join all of them to create the meta-analysis object.

3 Performing Meta-analysis

DExMA package contains the main gene expression meta-analysis methods:

- Meta-analysis based on effect size combination: Fixed Effects Model (FEM) and Random Effects Model (REM).
- Meta-analysis based on P-value combination: Fisher's method, Stouffer's method, Wilkerson's method (maxP) and Tippet's method (minP).

These methods can be applied directly, but it is advisable to apply some previous **DExMA** function to ensure the results are accurate.

3.1 Gene annotation and quality controls

Before performing the meta-analysis, all the genes must be in the same annotation and a quality control should be done in order to obtain reliable results [1]. **DExMA** provides some useful functions to help the user to do it.

3.1.1 Setting all the datasets in the same annotation

All genes must be in the same annotation in order to perform the meta-analysis successfully and avoid incorrect interpretations of the results. In cases where all datasets do not have the same gene ID, **DExMA** contains a function called *allSameID()* that allows to annotated all the datasets with the same gene ID. Specifically, the function allows to annotate those datasets that use the Official Gene Symbol, Entrez or Ensembl. The inputs of this function are:

- *objectMA*: the meta-analysis object of **DExMA** package. The result obtained by *createObjectMA()* function should be used.
- *initialIDs*: a character vector with the current annotation of each dataset. Use *availableIDs* function to see the annotations admitted.
- *finalID*: a character that indicates the final gene ID that all the studies will have.
- *organism*: a character that indicates the organism that the datasets belong to. Use *availableOrganism* function to see the organism admitted.

DExMA package

Here we include an example of how to use this function. In this example we have used "newObjectMA" that has been created before. The first two expression arrays are annotated in "entrez", the third expression matrix in "Official Gene Symbol" and the last one in "Official Gene Symbol" with some synonyms:

```
> rownames(newObjectMA$Study1$mExpres)[1:20]

[1] "100859927" "8086"      "8212"      "65985"     "729522"    "13"
[7] "344752"    "126767"    "343066"    "51166"     "79719"     "22848"
[13] "14"        "15"        "16"        "57505"     "80755"     "132949"
[19] "60496"     "10157"

> rownames(newObjectMA$Study2$mExpres)[1:20]

[1] "100859927" "8086"      "8212"      "65985"     "729522"    "13"
[7] "344752"    "126767"    "343066"    "51166"     "79719"     "22848"
[13] "14"        "15"        "16"        "57505"     "80755"     "132949"
[19] "60496"     "10157"

> rownames(newObjectMA$Study3$mExpres)[1:20]

[1] "AAA4"      "AAAS"      "AABT"      "AACS"      "AACSP1"    "AADAC"
[7] "AADACL2"   "AADACL3"   "AADACL4"   "AADAT"     "AAGAB"     "AAK1"
[13] "AAMP"      "AANAT"     "AARS1"     "AARS2"     "AARSD1"    "AASDH"
[19] "AASDHPPT" "AASS"

> rownames(newObjectMA$Study4$mExpres)[1:20]

[1] "AAA4"      "ADRACALIN" "AABT"      "ACSF1"     "AACSP1"    "AADAC"
[7] "AADACL2"   "AADACL3"   "AADACL4"   "AADAT"     "PPKP1"     "AAK1"
[13] "AAMP"      "AANAT"     "AARS1"     "AARS2"     "AARSD1"    "AASDH"
[19] "AASDHPPT" "AASS"
```

We can use *availableIDs* and *availableOrganism* in order to know how to write in the function the *initialIDs* vector and the organism:

```
> head(availableIDs)

[1] "Entrez"      "Ensembl"    "GeneSymbol"

> availableOrganism

[1] "Bos taurus"           "Caenorhabditis elegans"
[3] "Canis familiaris"     "Danio rerio"
[5] "Drosophila melanogaster" "Gallus gallus"
[7] "Homo sapiens"         "Mus musculus"
[9] "Rattus norvegicus"    "Arabidopsis thaliana"
[11] "Saccharomyces cerevisiae" "Escherichia coli"
```

We create the *InitialIDs* vector:

```
> annotations <- c("Entrez", "Entrez", "GeneSymbol", "GeneSymbol")
```

Then, we are ready to run the *allSameID()* function. We are going to annotated all the datasets in Official Gene Symbol (*finalID*="GeneSymbol"):

```
> newObjectMA <- allSameID(newObjectMA, initialIDs = annotations,
+                           finalID="GeneSymbol", organism = "Homo sapiens")
```

```

|
|
|
|=====| 25%
|
|=====| 50%
|
|=====| 75%
|
|=====| 100%

> rownames(newObjectMA$Study1$mExpres)[1:20]

[1] "AAA4"      "AAAS"      "AABT"      "AACS"      "AACSP1"    "AADAC"
[7] "AADACL2"   "AADACL3"   "AADACL4"   "AADAT"     "AAGAB"     "AAK1"
[13] "AAMP"      "AANAT"     "AARS1"     "AARS2"     "AARSD1"    "AASDH"
[19] "AASDHPPT" "AASS"

> rownames(newObjectMA$Study2$mExpres)[1:20]

[1] "AAA4"      "AAAS"      "AABT"      "AACS"      "AACSP1"    "AADAC"
[7] "AADACL2"   "AADACL3"   "AADACL4"   "AADAT"     "AAGAB"     "AAK1"
[13] "AAMP"      "AANAT"     "AARS1"     "AARS2"     "AARSD1"    "AASDH"
[19] "AASDHPPT" "AASS"

> rownames(newObjectMA$Study3$mExpres)[1:20]

[1] "AARS1"     "AATF"      "ABCC2"     "ABCD1P4"   "ABCD1P3"   "ABCD1P2"   "ACAD8"
[8] "AAVS1"     "ACAD9"     "ABT1"      "VSX1"      "ABHD5"     "AADAT"     "ABI3"
[15] "VRK3"      "VPS54"     "ABAT"      "VSI10"     "VRTN"      "VPS53"

> rownames(newObjectMA$Study4$mExpres)[1:20]

[1] "AARS1"     "AATF"      "ABCC2"     "ABCD1P4"   "ABCD1P3"   "ABCD1P2"   "ACAD8"
[8] "AAVS1"     "ACAD9"     "ABT1"      "VSX1"      "ABHD5"     "AADAT"     "ABI3"
[15] "VRK3"      "VPS54"     "ABAT"      "VSI10"     "VRTN"      "VPS53"

```

As it can be seen, all the studies are now annotated in Official Gene Symbol.

3.1.2 Logarithm transformation

To avoid problems with the returned fold-change by the meta-analysis, \log_2 should be applied to the gene expression values. We can make use of `dataLog()` function to check if each dataset expression values have the \log_2 applied already. If not, the function will make the transformation:

```

> newObjectMA <- dataLog(newObjectMA)
> head(newObjectMA[[1]][[1]])

      Sample1 Sample2 Sample3 Sample4
AAA4  5.439524 6.253319 2.926444 4.4304023
AAAS  5.769823 5.971453 1.831349 4.0466288
AABT  7.558708 5.957130 2.365252 3.4352889
AACS  6.070508 7.368602 2.971158 3.7151784
AACSP1 6.129288 5.774229 3.670696 3.9171749
AADAC  7.715065 7.516471 1.349453 0.3390772

```

3.1.3 Heterogeneity study

Some heterogeneity between studies may lead to some methods, such as the Fixed Effects Model, not providing reliable results. Therefore, it is advisable to carry out a study of heterogeneity to correctly choose the meta-analysis method [1]. The *heterogeneityTest()* function shows two ways of measuring heterogeneity.

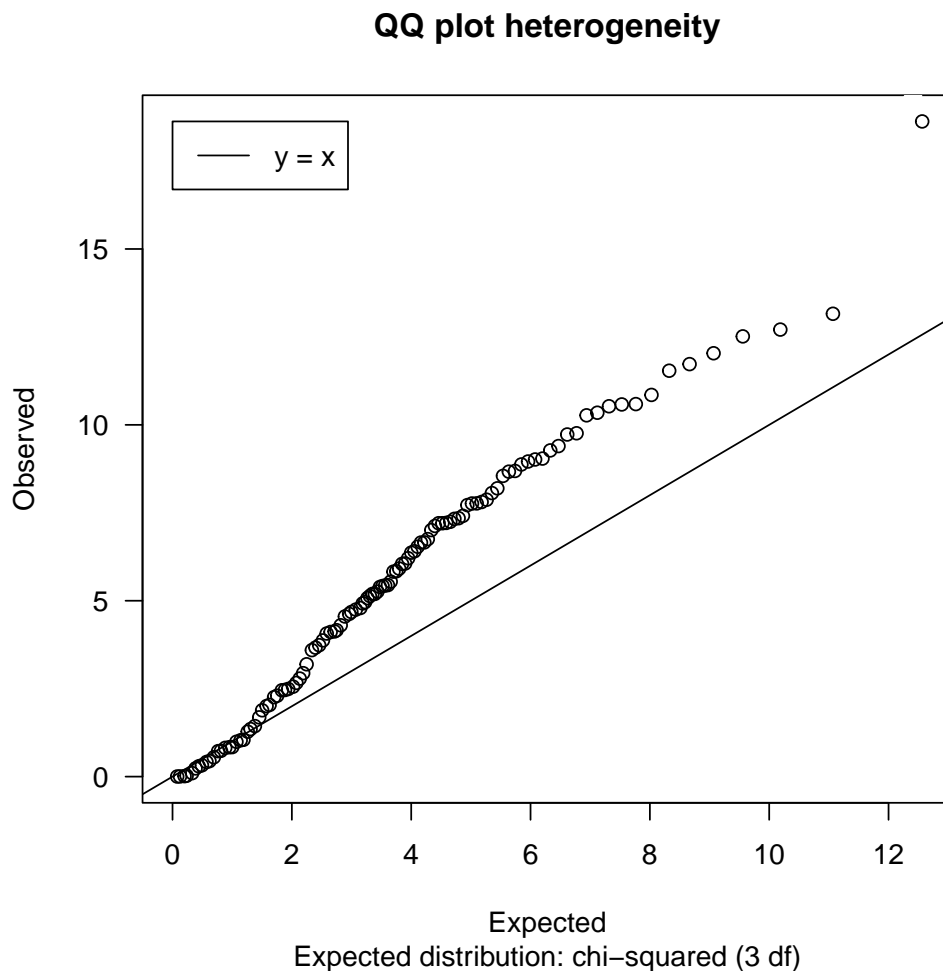
On the one hand, it returns a QQ-plot of the Cochran's test [4]. In this plot, if most of the values are close to the central line, that is, most of the Cochran's test values are close to the expected distribution (chi-squared distribution), it can be said that there is homogeneity. In the case that these values deviate greatly from the expected distribution, it must be assumed that there is heterogeneity.

On the other hand, I^2 measures the percentage of variation across studies due to heterogeneity [5]. In the case of gene expression data, an I^2 for each gene across datasets would have to be calculated. As in the case of many genes, it can be difficult to observe all the I^2 values obtained, the *heterogeneityTest()* function returns the quantiles of the different I^2 values calculated. I^2 values equal to 0 indicate homogeneity and values less than 0.25 are usually categorized as low heterogeneity [5]. Therefore, to assume homogeneity in the gene expression meta-analysis, almost all I^2 values must be 0 or at least less than 0.25.

In the example shown below, it is observed that in the QQ-plot of the Cochran's test, Q-values deviate considerably from the expected distribution and approximately 10% of the I^2 values are greater than 0.25, therefore homogeneity could not be assumed.

```
> heterogeneityTest(newObjectMA)

[1] "I^2 Quantiles"
      0%      5%     10%     15%     20%     25%     30%     35%
0.0000000 0.0000000 0.0000000 0.0000000 0.0000000 0.0000000 0.0000000 0.0000000
      40%     45%     50%     55%     60%     65%     70%     75%
0.0000000 0.0000000 0.1752591 0.2680708 0.3501887 0.4094427 0.4568317 0.5363021
      80%     85%     90%     95%    100%
0.5844494 0.6185599 0.6678108 0.7165413 0.8389331
```



3.2 Imputing missing genes expression values (optional)

Optionally, the *missGenesImput()* function allows to impute the expression values of the unmeasured genes. To perform this imputation, the function apply the knn method in the space of samples (*sampleKnn*) [2]. Once the function have been applied, it returns the same *objectMA* with all the datasets imputed:

```
> nrow(newObjectMA$Study1[[1]])
[1] 144
> nrow(newObjectMA$Study2[[1]])
[1] 144
> maObject_imput <- missGenesImput(maObject, k =7)
```

If we observe the number of genes in Study1, once the imputation is performed, we can observe that it goes from having 144 genes to 175.

```
> nrow(maObject_imput$Study1[[1]])
[1] 175

> nrow(maObject_imput$Study2[[1]])
[1] 175
```

3.3 Performing meta-analysis: *metaAnalysisDE()*

The *metaAnalysisDE()* function allows to perform a meta-analysis in only one step, needing only the meta-analysis object created previously.

This function has as input:

- **objectMA**: The meta-analysis object of DExMA package. The result obtained by *createObjectMA* function should be used.
- **typeMethod**: a character that indicates the method to be performed:
 - "FEM": Fixed Effects model.
 - "REM": Random Effects model.
 - "Fisher": Fisher's method (sum of logarithms of p-values)
 - "Stouffer": Stouffer's method (sum of z-scores)
 - "maxP": Wilkison's method (maximun of p-values)
 - "minP": Tippet's method (minimun of p-values)
- **missAllow**: a number between 0 and 1 that indicates the maximum proportion of missing values allows in a sample. If the sample has more proportion of missing values, the sample will be eliminated. In the other case, the missing values will be imputed by using the K-NN algorithm included in *impute* package [6]. In case the *objectMA* has been previously imputed, this element is not necessary.
- **proportionData**: a number between 0 and 1 that indicates the minimum proportion of datasets in which a gene must be contained to be included. In case the *objectMA* has been previously imputed, this element is not necessary.
- **Adjusted p-value** from which a gene is considered significant. Default 0.05. This value only takes you into account in the results generated in rank combination methods.

In the following example, we have applied a Random Effect model to the *DExMA* object ("newObjectMA") we have been working with so far. In addition we have allowed a 0.3 proportion of missing values in a sample and a gene must have been contained in at least the 50% of studies.

```
> resultsMA <- metaAnalysisDE(newObjectMA, typeMethod="REM",
+                             missAllow=0.3, proportionData=0.50)

[1] "Performing Random Effects Model"
```

The output of this function is a dataframe with the results of the meta-analysis where rows are the genes and columns are the different variables provided by the meta-analysis:

```
> head(resultsMA)
```

| | Com.ES | ES.var | Qval | Qpval | tau2 | Zval | Pval |
|--------|----------|-----------|-----------|------------|-----------|----------|--------------|
| AAA4 | 2.359381 | 0.7981372 | 6.0473427 | 0.10932948 | 1.5664238 | 2.640944 | 8.267527e-03 |
| AAAS | 2.281224 | 0.3370171 | 0.8210328 | 0.84442992 | 0.0000000 | 3.929542 | 8.510791e-05 |
| AABT | 2.493595 | 0.8816888 | 6.2077523 | 0.10192867 | 1.7585296 | 2.655635 | 7.915935e-03 |
| AACS | 2.313949 | 0.5884358 | 4.6107227 | 0.20262416 | 0.8221804 | 3.016506 | 2.557062e-03 |
| AACSP1 | 3.772840 | 0.5852004 | 2.0056235 | 0.57124034 | 0.0000000 | 4.931921 | 8.142482e-07 |
| AADAC | 2.697529 | 1.0092867 | 6.5349523 | 0.08829422 | 2.0628748 | 2.685090 | 7.251027e-03 |

| | FDR | propDataset |
|--------|--------------|-------------|
| AAA4 | 0.0401893697 | 1 |
| AAAS | 0.0018617356 | 1 |
| AABT | 0.0401893697 | 1 |
| AACS | 0.0178994351 | 1 |
| AACSP1 | 0.0001424934 | 1 |
| AADAC | 0.0396540550 | 1 |

Moreover, this function can also be applied to the imputed *objectMA*, but in this case, it is not necessary to indicate the *missAllow* and *proportionData* elements:

```
> resultsMA_imput <- metaAnalysisDE(maObject_imput, typeMethod="REM",)
[1] "Performing Random Effects Model"
> head(resultsMA_imput)
```

| | Com.ES | ES.var | Qval | Qpval | tau2 | Zval | Pval |
|--------|----------|-----------|-----------|------------|-----------|----------|--------------|
| AAA4 | 2.359381 | 0.7981372 | 6.0473427 | 0.10932948 | 1.5664238 | 2.640944 | 8.267527e-03 |
| AAAS | 2.281224 | 0.3370171 | 0.8210328 | 0.84442992 | 0.0000000 | 3.929542 | 8.510791e-05 |
| AABT | 2.493595 | 0.8816888 | 6.2077523 | 0.10192867 | 1.7585296 | 2.655635 | 7.915935e-03 |
| AACS | 2.313949 | 0.5884358 | 4.6107227 | 0.20262416 | 0.8221804 | 3.016506 | 2.557062e-03 |
| AACSP1 | 3.772840 | 0.5852004 | 2.0056235 | 0.57124034 | 0.0000000 | 4.931921 | 8.142482e-07 |
| AADAC | 2.697529 | 1.0092867 | 6.5349523 | 0.08829422 | 2.0628748 | 2.685090 | 7.251027e-03 |

| | FDR | propDataset |
|--------|--------------|-------------|
| AAA4 | 0.0498902520 | 1 |
| AAAS | 0.0021276978 | 1 |
| AABT | 0.0494745922 | 1 |
| AACS | 0.0279678674 | 1 |
| AACSP1 | 0.0001424934 | 1 |
| AADAC | 0.0469973985 | 1 |

The variables of the dataframe change from one type of meta-analysis to another. A more detailed explanation of these results will be addressed in the following sections.

3.3.1 Effects size combination results

The "FEM" and "REM" methods provide a dataframe with the variables:

- Com.ES: combined effect of the gene.
- ES.var: variance of the combined effect of the gene.
- Qval: total variance of the gene.
- Qpval:p-value for the total variance of the gene.

- tau2: between-study variance of the gene.
- zval: combined effect value for a standard normal. I can be use in order to find out if the gene is overexpressed (positive value) or underexpressed (negative value).
- Pval: P-value of the meta-analysis for the gene.
- FDR: P-value adjusted of the meta-analysis for the gene.
- Prop.dataset: Proportion of the datasets in which the gene is included.

```
> resultsES <- metaAnalysisDE(newObjectMA, typeMethod="REM", proportionData=0.5)
[1] "Performing Random Effects Model"
> head(resultsES)
```

| | Com.ES | ES.var | Qval | Qpval | tau2 | Zval | Pval |
|--------|--------------|-----------|-------------|------------|-----------|----------|--------------|
| AAA4 | 2.359381 | 0.7981372 | 6.0473427 | 0.10932948 | 1.5664238 | 2.640944 | 8.267527e-03 |
| AAAS | 2.281224 | 0.3370171 | 0.8210328 | 0.84442992 | 0.0000000 | 3.929542 | 8.510791e-05 |
| AABT | 2.493595 | 0.8816888 | 6.2077523 | 0.10192867 | 1.7585296 | 2.655635 | 7.915935e-03 |
| AACS | 2.313949 | 0.5884358 | 4.6107227 | 0.20262416 | 0.8221804 | 3.016506 | 2.557062e-03 |
| AACSP1 | 3.772840 | 0.5852004 | 2.0056235 | 0.57124034 | 0.0000000 | 4.931921 | 8.142482e-07 |
| AADAC | 2.697529 | 1.0092867 | 6.5349523 | 0.08829422 | 2.0628748 | 2.685090 | 7.251027e-03 |
| | FDR | | propDataset | | | | |
| AAA4 | 0.0401893697 | | 1 | | | | |
| AAAS | 0.0018617356 | | 1 | | | | |
| AABT | 0.0401893697 | | 1 | | | | |
| AACS | 0.0178994351 | | 1 | | | | |
| AACSP1 | 0.0001424934 | | 1 | | | | |
| AADAC | 0.0396540550 | | 1 | | | | |

3.3.2 P-value combination results

The "Fisher", "Stouffer", "minP" and "maxP" methods provide a dataframe with the following variables:

- Stat: Statistical calculated in the method
- Pval: P-value of the meta-analysis for the gene.
- FDR: P-value adjusted of the meta-analysis for the gene.
- AveFC: Average of log Fold-Change values for the gene used in order to find out if the gene is overexpressedd (positive value) or underexpressed (negative value).
- Prop.dataset: Proportion of the datasets in which the gene is included.

Here we present an example making use of "maxP" method:

```
> resultsPV <- metaAnalysisDE(newObjectMA, typeMethod="maxP", proportionData=0.5)
[1] "Performing MaxP's method"
> head(resultsPV)
```

| | Stat | Pval | FDR | AveFC | propDataset |
|------|------------|--------------|--------------|----------|-------------|
| AAA4 | 0.13754548 | 3.579194e-04 | 3.296626e-03 | 2.831428 | 1 |
| AAAS | 0.07824187 | 3.747632e-05 | 6.811215e-04 | 3.180932 | 1 |
| AABT | 0.27659986 | 5.853395e-03 | 2.768497e-02 | 2.935515 | 1 |

| | | | | | |
|--------|------------|--------------|--------------|----------|---|
| AACS | 0.11561455 | 1.786693e-04 | 1.737063e-03 | 3.060848 | 1 |
| AACSP1 | 0.01862270 | 1.202736e-07 | 1.052394e-05 | 2.953126 | 1 |
| AADAC | 0.22367842 | 2.503204e-03 | 1.510554e-02 | 3.629187 | 1 |

3.4 Visualization of the results: heatmap

Finally, we can represent in a heatmap the significant genes in order to observe how they are expressed in each of the studies. In `makeHeatma()` function we have to include both the object that has been used in the meta-analysis, the result of it and the applied method. In addition, this package offers three different scaling approaches (*scaling*) in order to compare properly the gene expression of the studies in the heatmap:

- "zscor": It calculates a z-score value for each gene, that is, the mean gene expression from each gene is subtracted from each gene expression value and then it is divided by the standard deviation.
- "swr": Scaling relative to reference dataset approach [7].
- "rscale": It uses the rescale function of the *scales* package to scale the gene expression [8].
- "none": no scaling approach is applied.

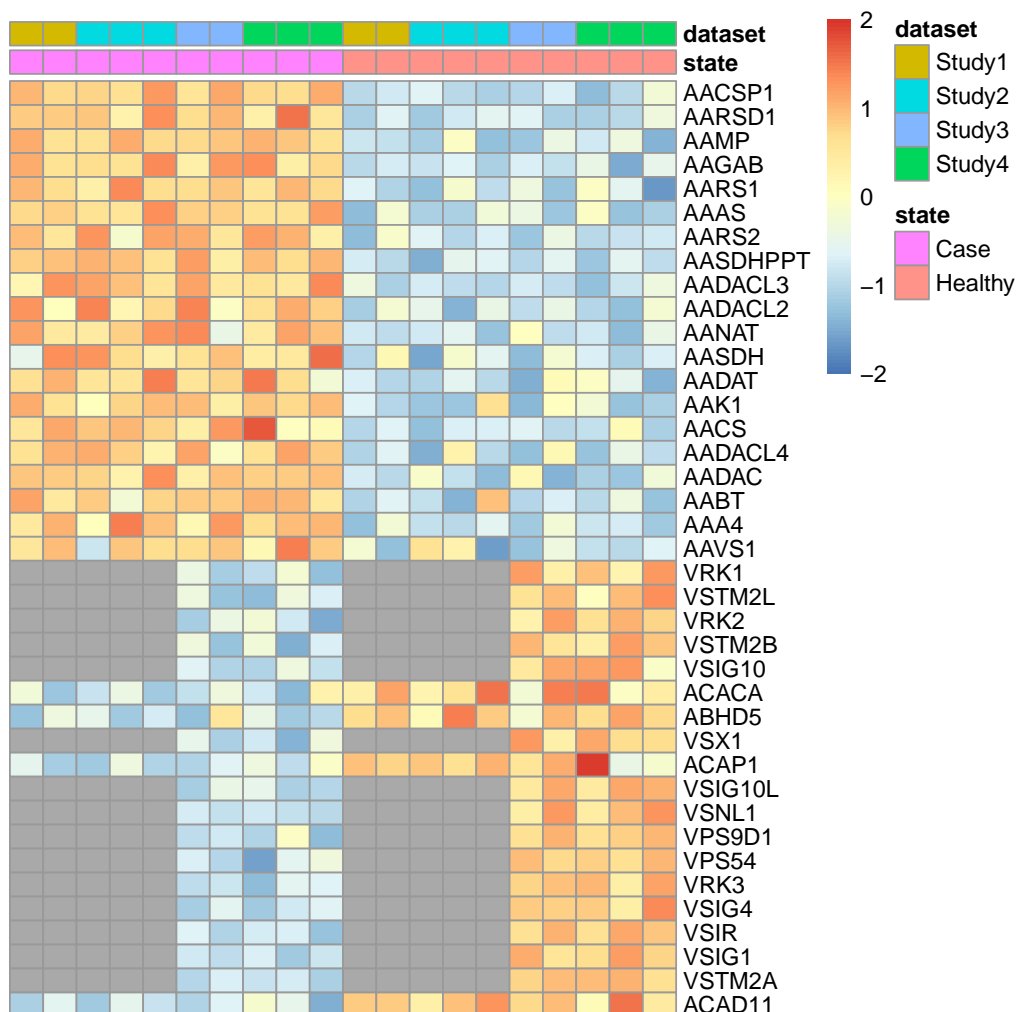
Moreover, in *regulation* argument, we can choose if we want to represent the overexpressed or underexpressed genes:

- "up": only up-expressed genes are represented.
- "down": only down-expressed genes are represented
- "all": up-expressed and down-expressed genes are represented.

We can choose the number of significant genes (*numSig*) that we want to be shown on the graph and the adjusted p-value from which a gene is considered as significant (*fdrSig*). In addition, the genes that are not presented in one sample are represented in gray.

Here we present an example of the heatmap which have been obtained from the result of applying a random effects model to the object "newObjectMA" and making use of a "zscor" scaling approach.

```
> makeHeatmap(objectMA=newObjectMA, resMA=resultsMA, scaling = "zscor",
+             regulation = "all", typeMethod="REM", numSig=40)
[1] "scaling using z-score..."
```



4 Additional information

DExMA provides some functions which may be useful for the user, although they are not essential to perform meta-analysis.

4.1 GEO microarray data download

In addition to using own user data, **DExMA** package allows to make use of public microarray data from the NCBI GEO public database [3]. For doing that, we can make use of `downloadGEOData()` function. This function uses internally `GEOquery` package in order to download some files at the same time. This function has an input a character vector (`GEOobjects`) with the GEO ID of the different datasets that we want to download and a character (`directory`) that indicates the directory where GSE Series Matrix files [9] are going to be stored.

```
> GEOobjects<- c("GSE4588", "GSE10325")
> dataGEO<-downloadGEOData(GEOobjects)
```

DExMA package

Once the download process is completed, we get a list of ExpressionSets. This list can be used as input of `createObjectMA()` function, although it is advisable to homogenize gene annotation, in case genes IDs are not Entrez, Official Gene Symbol or Ensembl.

4.2 Using RNA-Seq data

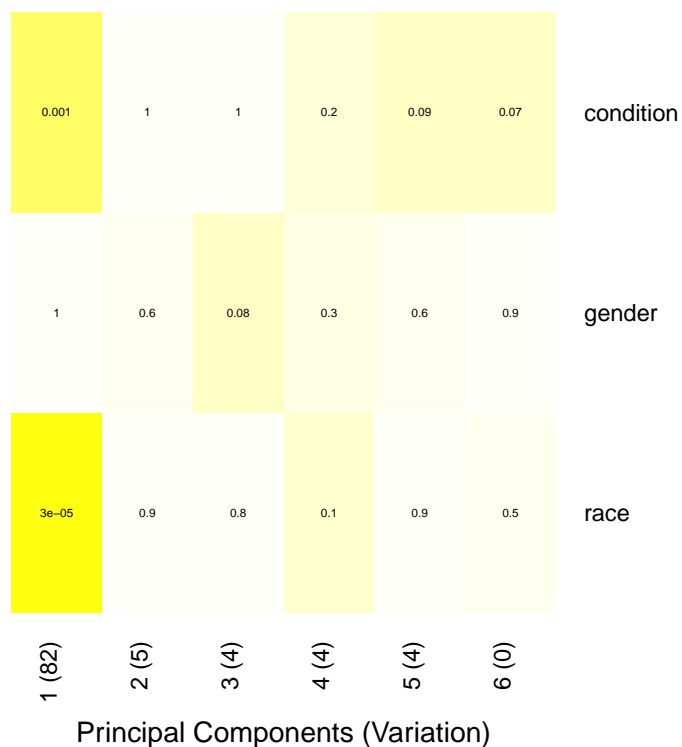
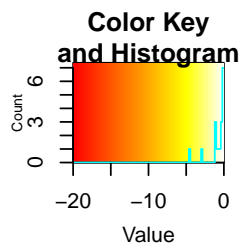
DExMA internally uses [limma](#) package in order to assess differential expression. Therefore, RNA-Seq data must be previously normalized by the user in order to be able to include correctly these data in the gene-expression meta-analysis. Since [limma](#) is used internally, we recommend to apply the steps described in the limma user's guide for the RNA-Seq data normalization [10], although the users can use the type of normalization they prefer.

4.3 Removing Batch Effects

Before the creation of the *objectMA*, a batch effect correction can be applied in order to reduce the effect of covariates that may be affecting to gene expression [1]. Firstly, with function `seeCov()`, which internally contains the `prince()` and `prince.plot()` functions of the [swamp](#) package [11], we can obtain a visualization of the p-values of each principal component associated with the categorical covariates. This allows to check which categorical variables are the ones that are most affecting the expression:

```
> seeCov(listMatrixEX$Study2, listPhenodatas$Study2)

Sample5 Sample6 Sample7 Sample8 Sample9 Sample10
      AA      AA      C      H      H      H
Levels: AA C H
```



The categorical variables that may be causing a batch effect can be corrected by using the `removeBatch()` function. The input of this function is the expression matrix and the phenodata. In addition, We also have to add a formula with the variables for which we want to correct the gene expression and the name of the variable that contains the cases and controls groups. Finally, if there is a covariate inside the formula that we want to give greater importance, we will have the option to indicate in the function (`mainCov()`). Here we show an example in which we have corrected the gene expression of the previous study by two of their covariates:

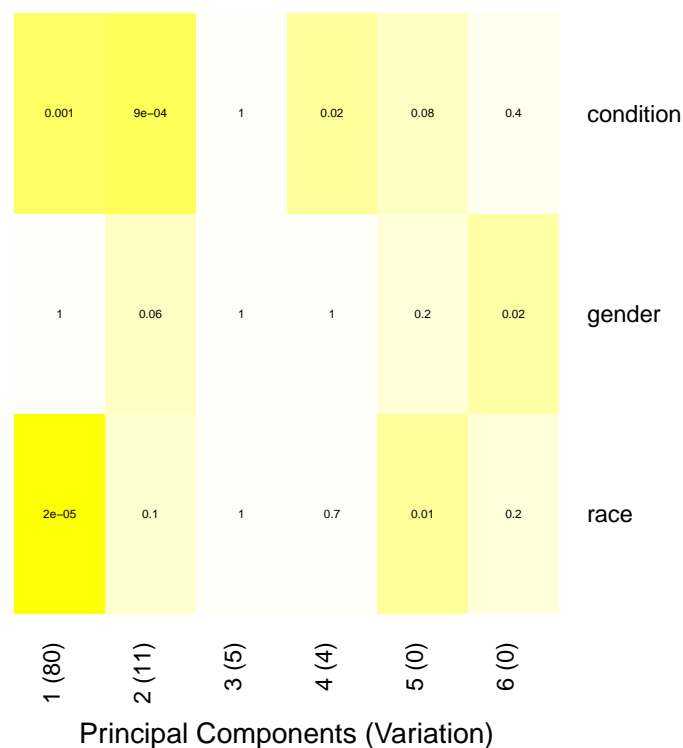
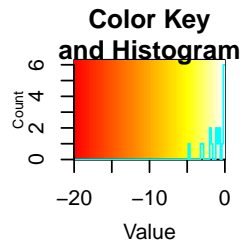
```
> listMatrixEX$Study2 <- batchRemove(listMatrixEX$Study2, listPhenodatas$Study2,
+                                     formula=~gender+race,
+                                     mainCov = "race", nameGroup="condition")

Coefficients not estimable: batch1 batch2 (Intercept) raceC raceH

> head(listMatrixEX$Study2)
```

| | Sample5 | Sample6 | Sample7 | Sample8 | Sample9 | Sample10 |
|-----------|----------|----------|----------|----------|----------|----------|
| 100859927 | 4.367690 | 6.648252 | 5.143704 | 2.820073 | 2.688687 | 2.688687 |
| 8086 | 5.937004 | 5.914434 | 5.831974 | 3.234439 | 3.228559 | 3.228559 |
| 8212 | 5.294553 | 4.092231 | 2.575022 | 3.368330 | 2.750491 | 2.750491 |

| | | | | | | |
|--------|----------|----------|----------|----------|----------|----------|
| 65985 | 5.685822 | 5.900363 | 5.544146 | 2.153294 | 3.166798 | 3.166798 |
| 729522 | 5.733054 | 5.553966 | 6.798683 | 3.275042 | 2.659061 | 2.659061 |
| 13 | 6.153159 | 5.435909 | 7.621115 | 4.860632 | 3.751608 | 3.751608 |



4.4 Calculating Effects size

The `calculateES()` function returns the effects size in each of the studies. Moreover, it calculates the variance of each of the effects and the proportion of datasets that contain the gene. The effects size are calculated by making use of the *Hedges' g estimator* [1].

```
> effects <- calculateES(newObjectMA)
> head(effects$ES)
```

| | Study1 | Study2 | Study3 | Study4 |
|------|----------|-----------|-----------|----------|
| AAA4 | 1.448942 | 2.5008385 | 1.1221262 | 6.642233 |
| AAAS | 1.506204 | 2.9643045 | 2.2593392 | 2.430575 |
| AABT | 2.288894 | 0.7656633 | 5.7273196 | 3.478798 |
| AACS | 2.579019 | 5.3315818 | 1.8723418 | 1.207949 |

```
AACSP1 5.705553 5.1083967 2.9008711 3.138604
AADAC 7.515582 2.1970319 0.8345501 3.776154

> head(effects$Var)

      Study1 Study2 Study3 Study4
AAA4 1.262429 1.187849 1.157396 4.3432711
AAAS 1.283581 1.398925 1.638077 1.1589745
AABT 1.654879 0.715520 5.100274 1.6751698
AACS 1.831418 3.035480 1.438208 0.7882618
AACSP1 5.069167 2.841310 2.051882 1.4875697
AADAC 8.060497 1.068912 1.087059 1.8549449
```

4.5 Calculating Individual P-values

Similar to the calculation of effects sizes, the individual p-values of each of the studies and the \log_2 fold change of each one can also be calculated by applying `pvalueIndAnalysis()`. P-value are obtained by assessing differential expression with *limma* package.

```
> pvalues <- pvalueIndAnalysis(newObjectMA)
> head(pvalues$p)

      Study1 Study2 Study3 Study4
AAA4 0.084586213 0.0124328637 0.137545477 0.0002392007
AAAS 0.078241874 0.0065691298 0.032924671 0.0139623839
AABT 0.031889060 0.2765998575 0.003696266 0.0035646909
AACS 0.024335148 0.0005986238 0.049748561 0.1156145550
AACSP1 0.003718102 0.0007176307 0.018622696 0.0053471186
AADAC 0.001900842 0.0197199607 0.223678418 0.0025641765

> head(pvalues$FC)

      Study1 Study2 Study3 Study4
AAA4 2.167998 2.654066 3.118286 3.385361
AAAS 2.931649 2.663952 4.321073 2.807056
AABT 3.857649 1.030831 3.259410 3.594168
AACS 3.376387 2.881147 4.252159 1.733698
AACSP1 2.157823 3.164180 3.998437 2.492064
AADAC 6.771502 2.282112 1.747278 3.715857
```

5 Session info

```
R version 4.2.0 RC (2022-04-19 r82224)
Platform: x86_64-pc-linux-gnu (64-bit)
Running under: Ubuntu 20.04.4 LTS
```

```
Matrix products: default
```

```
BLAS: /home/biocbuild/bbs-3.15-bioc/R/lib/libRblas.so
```

```
LAPACK: /home/biocbuild/bbs-3.15-bioc/R/lib/libRlapack.so
```

```
locale:
```

DExMA package

```
[1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
[3] LC_TIME=en_GB             LC_COLLATE=C
[5] LC_MONETARY=en_US.UTF-8   LC_MESSAGES=en_US.UTF-8
[7] LC_PAPER=en_US.UTF-8      LC_NAME=C
[9] LC_ADDRESS=C              LC_TELEPHONE=C
[11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
```

attached base packages:

```
[1] stats      graphics  grDevices  utils      datasets  methods    base
```

other attached packages:

```
[1] Biobase_2.56.0      BiocGenerics_0.42.0 DExMA_1.4.0
[4] DExMAdata_1.3.0
```

loaded via a namespace (and not attached):

```
[1] nlme_3.1-157          bitops_1.0-7          matrixStats_0.62.0
[4] bit64_4.0.5           RColorBrewer_1.1-3    httr_1.4.2
[7] GenomeInfoDb_1.32.0   tools_4.2.0           utf8_1.2.2
[10] R6_2.5.1              KernSmooth_2.23-20    DBI_1.1.2
[13] mgcv_1.8-40           colorspace_2.0-3      tidyselect_1.1.2
[16] bit_4.0.4             compiler_4.2.0         cli_3.3.0
[19] xml2_1.3.3            caTools_1.18.2        scales_1.2.0
[22] readr_2.1.2           genefilter_1.78.0     digest_0.6.29
[25] rmarkdown_2.14        GEOquery_2.64.0       XVector_0.36.0
[28] pkgconfig_2.0.3       htmltools_0.5.2       bnstruct_1.0.12
[31] fastmap_1.1.0         limma_3.52.0          rlang_1.0.2
[34] RSQLite_2.2.12        impute_1.70.0         farver_2.1.0
[37] generics_0.1.2        BiocParallel_1.30.0   gtools_3.9.2
[40] swamp_1.5.1           dplyr_1.0.8           RCurl_1.98-1.6
[43] magrittr_2.0.3        GenomeInfoDbData_1.2.8 Matrix_1.4-1
[46] Rcpp_1.0.8.3          munsell_0.5.0         S4Vectors_0.34.0
[49] fansi_1.0.3           lifecycle_1.0.1       yaml_2.3.5
[52] edgeR_3.38.0          MASS_7.3-57           zlibbioc_1.42.0
[55] plyr_1.8.7            gplots_3.1.3          grid_4.2.0
[58] blob_1.2.3            parallel_4.2.0        snpStats_1.46.0
[61] crayon_1.5.1          lattice_0.20-45       Biostrings_2.64.0
[64] splines_4.2.0         annotate_1.74.0        hms_1.1.1
[67] KEGGREST_1.36.0       locfit_1.5-9.5        knitr_1.38
[70] pillar_1.7.0          igraph_1.3.1          stats4_4.2.0
[73] XML_3.99-0.9          glue_1.6.2            evaluate_0.15
[76] data.table_1.14.2     BiocManager_1.30.17   tzdb_0.3.0
[79] png_0.1-7            vctrs_0.4.1           tidyr_1.2.0
[82] gtable_0.3.0          purrr_0.3.4           assertthat_0.2.1
[85] amap_0.8-18           cachem_1.0.6          xfun_0.30
[88] xtable_1.8-4          survival_3.3-1        tibble_3.1.6
[91] pheatmap_1.0.12       AnnotationDbi_1.58.0  memoise_2.0.1
[94] IRanges_2.30.0        sva_3.44.0            ellipsis_0.3.2
[97] BiocStyle_2.24.0
```

References

- [1] Toro-Domínguez D., Villatoro-García J.A., Martorell-Marugán J., and et al. A survey of gene expression meta-analysis: methods and applications. *Briefings in Bioinformatics*, pages 1–12, 2020. doi:<https://doi.org/10.1093/bib/bbaa019>.
- [2] Christopher A. Mancuso, Jacob L. Canfield, Deepak Singla, and Arjun Krishnan. A flexible, interpretable, and accurate approach for imputing the expression of unmeasured genes. *Nucleic Acids Research*, 48(21):e125–e125, 2020. URL: <https://doi.org/10.1093/nar/gkaa881>.
- [3] Barret T., Wilhite S., Ledoux P., and et al. Ncbi geo: archive for functional genomics data sets–update. *Nucleic Acids Research*, pages 991–995, 2020. doi:<https://doi.org/10.1093/nar/gks1193>.
- [4] Higgins J. and Thompson S. Quantifying heterogeneity in a meta-analysis. *Statistics in Medicine*, pages 1539–1558, 2002. doi:[10.1002/sim.1186](https://doi.org/10.1002/sim.1186).
- [5] Higgins J., Thompson S., Deeks J., and Altman D. Measuring inconsistency in meta-analyses. *BMJ*, pages 557–560, 2003. doi:[10.1136/bmj.327.7414.557](https://doi.org/10.1136/bmj.327.7414.557).
- [6] Hastie T., Tibshirani R., Narasimhan B., and Chu G. impute: Imputation for microarray data. 2019.
- [7] Lazar C., Meganck S., Taminau J., and et al. Batch effect removal methods for microarray gene expression data integration: a survey. *Briefings in Bioinformatics*, pages 469–490, 2013. doi:[10.1093/bib/bbs037](https://doi.org/10.1093/bib/bbs037).
- [8] Wickham H. and Siedel D. Scale functions for visualization. 2020.
- [9] Davis S. and Meltzer P.S. Geoquery: a bridge between the gene expression omnibus (geo) and bioconductor. *Bioinformatics*, pages 1846–1847, 2007. doi:[10.1093/bioinformatics/btm254](https://doi.org/10.1093/bioinformatics/btm254).
- [10] Smyth G. K., Ritchie M., Thorne N., Yifang Hu, and et al. Linear models for microarray and rna-seq data user's guide. 2020.
- [11] Lauss M. Visualization, analysis and adjustment of high-dimensional data in respect to sample annotations. 2019.