

GeneGeneInteR vignette

Introduction

Mathieu Emily and Magalie Houée-Bigot

April 7, 2020

This vignette gives an overview of the main functions and tools proposed in the package **GeneGeneInteR**. Readers that would like to have more technical details about the statistical methods implemented in the package are encouraged to read the vignette **GenePair**.

1 Introduction

The package **GeneGeneInteR** aims at providing a collection of statistical methods to search for interaction between genes in case-control association studies. These methods are dedicated to the analysis of biallelic SNPs (Single Nucleotide Polymorphism) genotype data. This vignette describes the complete analysis pipeline of a set of $k > 2$ genes, going from data importation to results visualization via data manipulation and statistical analysis.

In the remainder of this vignette, we illustrate our pipeline through the analysis of a case-control dataset publicly available in the NCBI repository GSE39428 series [Chang et al., 2013]. The dataset contains the genotypes of 312 SNPs from 17 genes in a total of 429 patients (266 individuals affected by Rheumatoid Arthritis and 163 Health controls) and is attached to our package **GeneGeneInteR** as an external file in the ped format.

2 Data importation and manipulation

2.1 Importation of genotype data

At first, the path for the files containing genotype data and information regarding the SNP set are loaded:

```
> ped <- system.file("extdata/example.ped", package="GeneGeneInteR")
> info <- system.file("extdata/example.info", package="GeneGeneInteR")
> ## Information about position of the snps
> posi <- system.file("extdata/example.txt", package="GeneGeneInteR")
```

The importation is performed with the `importFile` function:

```
> data <- importFile(file=ped, snps=info, pos=posi, pos.sep="\t")
```

The object `data` is a list of 3 elements: `status`, `snpX` (a `Snpmatrix` object) and `genes.info` (a `data.frame`). The `status` is available only if the imported format is `ped`.

```
> summary(data)
```

| | Length | Class | Mode |
|-------------------------|--------|-------------------------|----------------------|
| <code>status</code> | 429 | <code>factor</code> | <code>numeric</code> |
| <code>snpX</code> | 133848 | <code>Snpmatrix</code> | <code>raw</code> |
| <code>genes.info</code> | 4 | <code>data.frame</code> | <code>list</code> |

We can check that the `snpX` object contains the genotype of 429 individuals for 312 SNPs.

```
> data$snpX

A SnpMatrix with 429 rows and 312 columns
Row names: H97 ... RA345
Col names: rs1002788 ... rs9502656
```

The `genes.info` is a `data.frame` with exactly four columns that are named as follows: `Chromosome`, `Genenames`, `SNPnames` and `Position`.

```
> summary(data$genes.info)
```

| | Chromosome | Genenames | SNPnames | Position |
|---------|------------|--------------|----------------|--------------------|
| Min. | : 1.000 | PCSK6 : 74 | rs1002788 : 1 | Min. : 7881078 |
| 1st Qu. | : 6.000 | TXNDC5 : 69 | rs1005753 : 1 | 1st Qu. : 11712674 |
| Median | : 8.000 | DNAH9 : 41 | rs1006273 : 1 | Median : 47863803 |
| Mean | : 9.962 | CA1 : 38 | rs10152164 : 1 | Mean : 52968484 |
| 3rd Qu. | : 15.000 | VDR : 19 | rs10184179 : 1 | 3rd Qu. : 97191582 |
| Max. | : 17.000 | Gc : 12 | rs1032551 : 1 | Max. : 123157722 |
| | | (Other) : 59 | (Other) : 306 | |

2.2 Phenotype importation

Similar to functions introduced to analyze a single pair of genes (see vignette “Statistical analysis of the interaction between a pair of genes.”), the case-control status is stored in a numeric or a factor vector with exactly two distinct values. If the phenotype is saved in a separate file in table form, it can thus be imported simply by using the `read.table` such as for example:

```
> Y <- read.table(system.file("/extdata/response.txt", package="GeneGeneInter"), sep=";")
```

If the case-control status is provided in the ped file, it can be uploaded as follows:

```
> Y <- data$status
```

2.3 Data filtering

Before performing the statistical analysis, it is very common to remove some SNPs in order to improve the quality of the data. Such a cleaning step can be performed in our `GeneGeneInter` package by using the function `snpMatrixScour`. `snpMatrixScour` aims at modifying a `SnpMatrix` object by removing SNPs that does not meet criteria regarding the Minor Allele Frequency (MAF), deviation to Hardy-Weinberg Equilibrium (HWE) and the proportion of missing values. In the following example, SNPs with MAF lower than 0.05 or SNPs with p-value for HWE lower than 0.001 or SNPs with a call rate lower than 0.9 are removed from the object `data`.

```
> data <- snpMatrixScour(data$snpX, genes.info=data$genes.info, min.maf=0.05
+ , min.eq=1e-3, call.rate=0.9)
```

The following lines show that the dataset now contains only 209 SNPs, meaning that 103 SNPs have been filtered out.

```
> data$snpX

A SnpMatrix with 429 rows and 209 columns
Row names: H97 ... RA345
Col names: rs10510123 ... rs4328262
```

Since the use of stringent filters could lead to the elimination of all SNPs within a gene, care has to be taken during the filtering step. However, in such a situation the gene without SNP is removed from the dataset and a warning message is provided for the user.

In other situations, the user might be interested in performing the analysis on a predefined subset of SNPs. For that purpose, the `selectSnps` function provides three options to extract of collection of SNPs by specifying the argument `select` that should be one of the following:

- a numeric vector with only the column number in the `snpMatrix` (or row number for `genes.info`) of each selected SNP. The following line allow the extraction of the 10 first SNPs:

```
> selec <- selectSnps(data$snpX, data$genes.info, select=1:10)
```

- a character vector with the names of each selected SNP or each selected gene. The following example is used to extract genes DNAH9 and TXNDC5:

```
> selec <- selectSnps(data$snpX, data$genes.info, c("DNAH9", "TXNDC5"))
```

- a character vector which elements are position bounds of genes. Each element of the vector is either of the form “begin:end”, or “chr:begin:end” if you have to precise the chromosome of the gene. The following code allow to select SNPs from position 101342000 to 101490000 on chromosome 15:

```
> selec <- selectSnps(data$snpX, data$genes.info, c("15:101342000:101490000"))
```

2.4 Imputation

Since our pipeline of analysis does not handle with missing values, SNPs filtering as well as SNPs selection can help removing missing data. This can be done easily by applying the `snpMatrixScour` with `call.rate=1` argument. However, in that case, SNPs with an acceptable call rate are also removed and the lost information is likely to be critical. Genotype imputation is then commonly performed to keep most of the informative SNPs in the dataset. Since our genotype data are stored into `SnpMatrix` object, we implement the `imputeSnpMatrix` function that wraps `snp.imputation` and `impute.snps` functions from `snpStats` package. Our `imputeSnpMatrix` function mimics a Leave-One-Out process where missing SNP are imputed for an individual based on a model trained on all other individuals.

In our example, the following lines show that after the filtering step, 844 missing values still remain in the dataset.

```
> sum(is.na(data$snpX))
```

```
[1] 844
```

To impute those missing values, we simply used our `imputeSnpMatrix` function as follows:

```
> data <- imputeSnpMatrix(data$snpX, data$genes.info)
```

```
|-----| 100%
```

A simple check of the dataset show that all missing values have been imputed:

```
> sum(is.na(data$snpX))
```

```
[1] 0
```

When the amount of missing values is so important that `snp.imputation` is not able to find a rule of imputation, some missing values may remain. In that case, the user can specify the action to be done thanks to the `om.rem` arguments:

- `om.rem="none"`: leave the dataset as it is,
- `om.rem="SNP"`: remove all SNPs with remaining missing values,
- `om.rem="ind"`: remove all individuals with remaining missing values.

It is noteworthy that removing all SNPs is often more parsimonious than removing individuals and allows to get a dataset without any missing values with minimum information-loss.

Although, function `snp.imputation` can calculate accurate rules for imputation, we encouraged the user to first input missing genotype with an external software (such as IMPUTE2 [Howie et al., 2009]) prior to the importation step.

3 Statistical analysis

The statistical analysis of a set of genes, as implemented in the `GGI` function, consists in performing all possible pairwise tests between two genes. Pairwise tests are conducted by using the `method` argument with one of the ten methods detailed in the vignette “Statistical analysis of the interaction between a pair of genes”. The `GGI` function takes two further mandatory arguments: `Y` the vector of case-control status and `snpX`, a `Snpmatrix` object that store the genotypes for all SNPs. It is assumed that SNPs within the same gene are consecutive in the `snpX` argument. Furthermore, gene information, such as gene ordering and the number of SNPs within each gene, has to be provided either in the `genes.length` or in the `gene.info` argument.

The following line allow the computation of all pairwise tests between the 17 genes of our example dataset with the PCA-based method.

```
> GGI.res <- GGI(Y=Y, snpX=data$snpX, genes.info=data$genes.info,method="PCA")
```

The output of the `GGI` function is a an object of class `GGInetwork`.

```
> class(GGI.res)
```

```
[1] "GGInetwork"
```

The class `GGInetwork` is an S3 class based on a list of 4 elements `statistic`, `p.value`, `method` and `parameter`. When `method="PCA"`, a fifth element, called `df`, is added to the `GGInetwork`.

```
> names(GGI.res)
```

```
[1] "statistic" "p.value" "df" "method" "parameter"
```

Elements `statistic`, `p.value` and `statistic`, `df` are squared matrices with M rows and M columns where M is the number of genes in the dataset. The general terms of each matrices are respectively the statistic, the p-value and the degrees of freedom of the Likelihood Ratio Test. The element `method` is the name of the method used to perform the pairwise interaction tests. Finally, the element `parameter` is a list of the parameters used to perform the pairwise interaction tests.

As example, the `GGI.res` object generated in the previous example is a list of 5 elements. Each cell of the output `p.value` matrix is the p-value of the corresponding pairwise test. The pairwise p-values obtained for the 4 first genes (bub3, CA1, CDSN, DNAH9) of our dataset can be observed as follows:

```
> round(GGI.res$p.value[1:4,1:4],digits=4)
```

| | bub3 | CA1 | CDSN | DNAH9 |
|-------|--------|--------|--------|--------|
| bub3 | 0.0000 | 0.1684 | 0.3179 | 0.1851 |
| CA1 | 0.1684 | 0.0000 | 0.0697 | 0.0000 |
| CDSN | 0.3179 | 0.0697 | 0.0000 | 0.4539 |
| DNAH9 | 0.1851 | 0.0000 | 0.4539 | 0.0000 |

Significant results can be summarized using the S3 method `summary` for class `GGInetwork`. The method `summary` prints the pairs of genes with a interaction p-value lower than 0.05 after (1) no correction (2) a bonferroni correction and (3) Benjamini & Hochberg correction for multiple testing.

```
> summary(GGI.res)
```

Gene-gene interaction network of 17 genes performed with:
Principal Component Analysis
Significant interaction with no correction at the level of 0.05

```
-----
      Gene1   Gene2  Uncorrected p-value
1  TXNDC5     VDR      8.9e-10
2  DNAH9    TXNDC5      1.9e-09
:
:
Significant interaction with a bonferroni correction at the level of 0.05
-----
      Gene1   Gene2  bonferroni p-value
1  TXNDC5     VDR      1.2e-07
2  DNAH9    TXNDC5      2.6e-07
:
:
Significant interaction with a Benjamini & Hochberg correction at the level of 0.05
-----
      Gene1   Gene2  BH p-value
1  TXNDC5     VDR      1.2e-07
2  DNAH9    TXNDC5      1.3e-07
```

4 Visualization

The visualization of the results can be performed with the S3 method `plot` for class `GGInetwork`. Given a `GGInetwork` object obtained from the analysis of M genes with our `GGI` function, results can be visualized through two types of representation: an heatmap-like visualization with the `method="heatmap"` argument and a network-like representation with the `method="network"` argument.

4.1 Heatmap-like visualization

The `plot` method can be simply used with the `GGInetwork` object as the single input argument. Figure 1 (a) and (b) show the graphical representation where all pairwise interactions are plotted (Fig. 1 (a)) or only the interaction between the 3 genes CA1, Gc and PADI1 with the argument `genes` (Fig. 1 (b)).

(a) `> plot(GGI.res)` (b) `> plot(GGI.res, genes=c("CA1", "Gc", "PADI1"))`

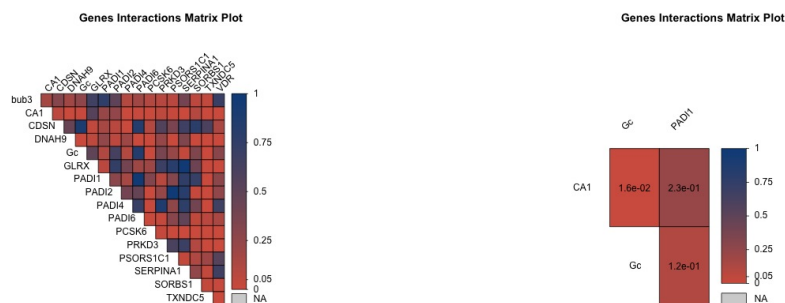


Figure 1: Default output of the function `GGI.plot` considering (a) the whole set of genes or (b) a restricted set of 3 genes.

When the number of genes is below 15, p-values and names are drawn to make matrix reading easier (see Figure 1(b)). However, when the number of genes is larger than 15, p-values are not drawn and gene names are kept while if the number of genes is larger than 25, none of the p-values or the gene names are displayed (see Figure 1(a)). In that case, the default behavior of the function is to start an interactive process where user can click on a cell of interest to open a tooltip displaying which genes are involved in the selected interaction and the p-value of the interaction test. Tooltips can be closed if user clicks anywhere else than on a cell. This process stops when the user presses the escape button (or terminates the locator procedure in general) or when the user clicks on any place other than a cell when no tooltip window is open.

Several arguments can further be specified to customize the output graphics such as colors (arguments `col` and `NA.col`), width of the bar for colors (argument `colbar.width`), and titles (argument `title`). User can also decide whether p-values (argument `draw.pvals`) and gene names (argument `draw.names`) should be drawn and is allowed to disable the interactivity of plot (argument `interact`). To further improve plot clarity and hence allowing a better interpretation of the results, (1) genes can be ordered according to a hierarchical clustering (argument `hclust.order`), (2) p-values can be reported in $-\log_{10}$ scale (argument `use.log`) and (3) a threshold can be applied to the p-values in order to distinguish between significant and non-significant interactions (argument `threshold`).

Figures 2 and 3 provide two plots resulting from different sets of arguments passed to the `plot` function.

```
> plot(GGI.res,col=c("black","cyan","white"),colbar.width=0.25,title="Interaction
between 17 genes",hclust.order=TRUE,use.log=TRUE,threshold=NULL,NA.col
="#D3D3D3",draw.pvals=FALSE,draw.names=TRUE,interact=FALSE)
```

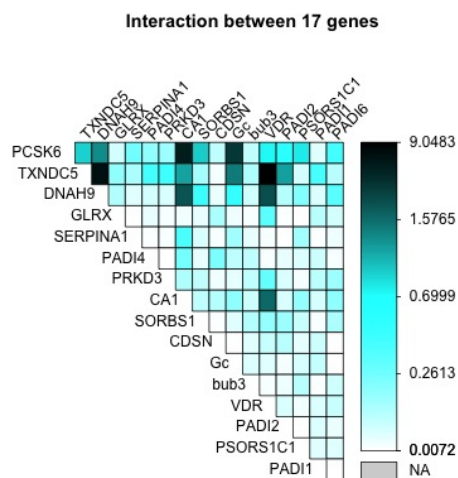


Figure 2: Example of the use of `plot` arguments when no threshold is applied to the p-values.

```
> plot(GGI.res,col=c("black","cyan","white"),colbar.width=0.05,title="Interaction
between 17 genes",hclust.order=TRUE,use.log=FALSE,threshold=0.05,NA.col
="#D3D3D3",draw.pvals=FALSE,draw.names=TRUE,interact=FALSE)
```

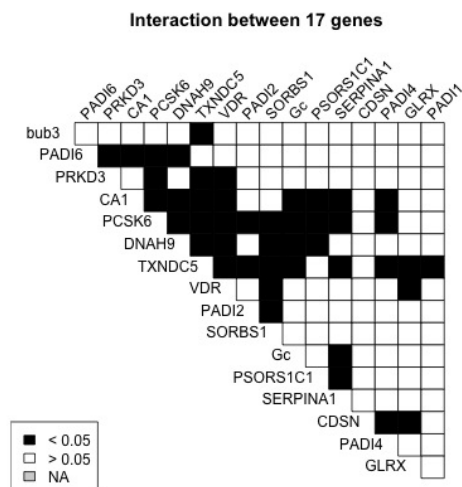


Figure 3: Example of the use of `plot` arguments when a threshold of 0.05 is applied to the p-values.

4.2 Network-like visualization

The `plot` function with `method="network"` aims at drawing a graph between genes where two genes are adjacent if the p-values between these two genes is below a given threshold (argument `threshold` with a default value equal to 0.05). The display of the network is performed by utilizing the `graph_from_data_frame` from the `igraph` R package [Csardi and Nepusz, 2006].

Two additional arguments can be used to customize the network. First, user can focus on a specific subset of genes with the argument `genes` and secondly, gene(s) not linked to other genes can be removed from the graph with argument `plot.nointer`.

Figure 4 displays the default network obtained with all genes. In figure 5, a subset of only 12 genes have been selected to be the vertices of the graph. However, genes `bub3` and `PADI1` does not have a pvalue below the threshold of 0.05 with any of the other selected genes. Since the argument `plot.nointer` is set to `TRUE`, the two genes `bub3` and `PADI1` are not drawn in the resulting network.

```
>set.seed(1234)
>plot(GGI.res,method="network")
```

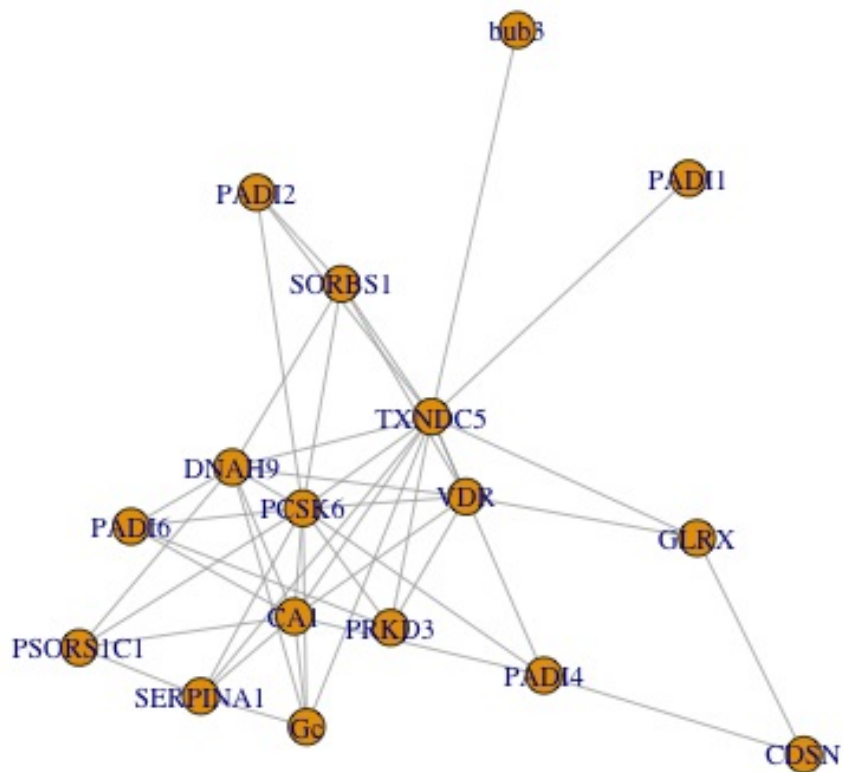


Figure 4: Default output of `draw.network`

```

> set.seed(1234)
> plot(GGI.res,method="network",genes=c("bub3","CDSN","Gc","GLRX","PADI1","PADI2","PADI4",
    "PADI6","PRKD3","PSORS1C1","SERPINA1","SORBS1"),
    threshold=0.05,plot.nointer=FALSE)

```

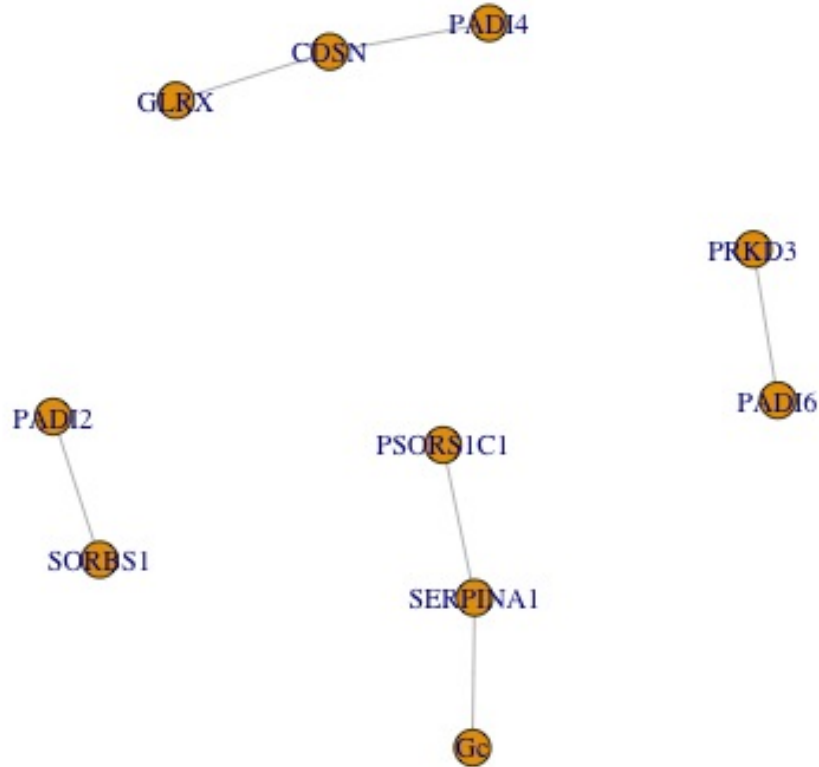


Figure 5: Output of `draw.network` with a customized set of arguments.

References

- [Chang et al., 2013] Chang, X., Xu, B., Wang, L., Wang, Y., Wang, Y., and Yan, S. (2013). Investigating a pathogenic role for *txndc5* in tumors. *International Journal of Oncology*, 43(43):1871–1884.
- [Csardi and Nepusz, 2006] Csardi, G. and Nepusz, T. (2006). The igraph software package for complex network research. *InterJournal*, Complex Systems:1695.
- [Howie et al., 2009] Howie, B. N., Donnelly, P., and Marchini, J. (2009). A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLoS Genetics*, 5(6):e1000529.