

Bioconductor's SPIA package

Adi L. Tarca^{1,2,3}, Purvesh Khatri¹ and Sorin Draghici¹

March 23, 2011

¹Department of Computer Science, Wayne State University

²Bioinformatics and Computational Biology Unit of the NIH Perinatology Research Branch

³Center for Molecular Medicine and Genetics, Wayne State University

1 Overview

This package implements the Signaling Pathway Impact Analysis (SPIA) algorithm described in Tarca et al. (2009), Khatri et al. (2007) and Draghici et al. (2007). SPIA uses the information from a set of differentially expressed genes and their fold changes, as well as pathways topology in order to assess the significance of the pathways in the condition under the study. The current version of SPIA algorithm uses KEGG signaling pathway data for several popular organisms. SPIA ready KEGG pathway data for more organisms is available at

<http://bioinformaticsprb.med.wayne.edu/SPIA/>.

The pathways included for each organism are those i) containing at least one relation between genes/proteins considered by SPIA, and ii) having no reactions.

The KEGG data that was preprocessed for SPIA analysis was downloaded from KEGG's ftp repository on: 03/21/2011. For a list of changes in SPIA compared to previous versions see the last section in this document.

2 Pathway analysis with SPIA package

This document provides basic introduction on how to use the SPIA package. For extended description of the methods used by this package please consult these references: Tarca et al. (2009); Khatri et al. (2007); Draghici et al. (2007).

We demonstrate the functionality of this package using a colorectal cancer dataset obtained using Affymetrix GeneChip technology and available through GEO (GSE4107). The experiment contains 10 normal samples and 12 colorectal cancer samples and is described by Hong et al. (2007). RMA preprocessing of the raw data was performed using the `affy` package, and a two group moderated t-test was applied using the `limma` package. The data frame obtained as an end result from the function `topTable` in `limma` is used as starting point for preparing the input data for SPIA. This data frame called `top` was made available in the `colorectalcancer` dataset included in the SPIA package:

```
> library(SPIA)
> data(colorectalancer)
> options(digits = 3)
> head(top)
```

	ID	logFC	AveExpr	t	P.Value	adj.P.Val	B	ENTREZ
10738	201289_at	5.96	6.23	23.9	1.79e-17	9.78e-13	25.4	3491
18604	209189_at	5.14	7.49	17.4	1.56e-14	2.84e-10	21.0	2353
11143	201694_s_at	4.15	7.04	16.5	5.15e-14	7.04e-10	20.1	1958
10490	201041_s_at	2.43	9.59	14.1	1.29e-12	1.41e-08	17.7	1843
10913	201464_x_at	1.53	8.22	11.0	1.69e-10	1.15e-06	13.6	3725
11463	202014_at	1.43	5.33	10.5	4.27e-10	2.42e-06	12.8	23645

For SPIA to work, we need a vector with log2 fold changes between the two groups for all the genes considered to be differentially expressed. The names of this vector must be Entrez gene IDs. The following lines will add one additional column in the `top` data frame annotating each affymetrix probeset to an Entrez ID. Since there may be several probesets for the same Entrez ID, there are two easy ways to obtain one log fold change per gene. The first option is to use the fold change of the most significant probeset for each gene, while the second option is to average the log fold-changes of all probesets of the same gene. In the example below we used the former approach. The genes in this example are called differentially expressed provided that their FDR adjusted p-values (q-values) are less than 0.05. The following lines start with the `top` data frame and produce two vectors that are required as input by `spia` function:

```
> library(hgu133plus2.db)
> x <- hgu133plus2ENTREZID
> top$ENTREZ <- unlist(as.list(x[top$ID]))
> top <- top[!is.na(top$ENTREZ), ]
> top <- top[!duplicated(top$ENTREZ), ]
> tg1 <- top[top$adj.P.Val < 0.1, ]
> DE_Colorectal = tg1$logFC
> names(DE_Colorectal) <- as.vector(tg1$ENTREZ)
> ALL_Colorectal = top$ENTREZ
```

The `DE_Colorectal` is a vector containing the log2 fold changes of the genes found to be differentially expressed between cancer and normal samples, and `ALL_Colorectal` is a vector with the Entrez IDs of all genes profiled on the microarray. The names of the `DE_Colorectal` are the Entrez gene IDs corresponding to the computed log fold-changes.

```
> DE_Colorectal[1:10]
```

3491	2353	1958	1843	3725	23645	9510	84869	7432	1490
5.96	5.14	4.15	2.43	1.53	1.43	3.94	-1.15	4.72	3.45

```
> ALL_Colorectal[1:10]
```

[1]	"3491"	"2353"	"1958"	"1843"	"3725"	"23645"	"9510"	"84869"	"7432"
[10]	"1490"								

The SPIA algorithm takes as input the two vectors above and produces a table of pathways ranked from the most to the least significant. This can be achieved by calling the `spia` function as follows:

```
> res = spia(de = DE_Colorectal, all = ALL_Colorectal, organism = "hsa",
+           nB = 2000, plots = FALSE, beta = NULL, combine = "fisher",
+           verbose = FALSE)
> res$Name = substr(res$Name, 1, 10)
> res[1:20, -12]
```

	Name	ID	pSize	NDE	pNDE	tA	pPERT	pG	pGFdr
1	Focal adhe	04510	177	88	7.26e-08	102.61	0.000005	1.08e-11	1.20e-09
2	Alzheimer'	05010	151	85	4.50e-11	-5.67	0.284000	3.34e-10	1.13e-08
3	ECM-recept	04512	74	42	2.63e-06	26.07	0.000005	3.43e-10	1.13e-08
4	Parkinson'	05012	110	65	5.42e-10	-13.49	0.029000	4.07e-10	1.13e-08
5	Pathways i	05200	304	128	1.23e-05	67.18	0.005000	1.08e-06	2.41e-05
6	PPAR signa	03320	65	37	9.34e-06	-3.11	0.046000	6.73e-06	1.24e-04
7	Huntington	05016	169	80	3.37e-06	-3.07	0.248000	1.25e-05	1.99e-04
8	Axon guida	04360	119	59	1.09e-05	11.54	0.166000	2.57e-05	3.56e-04
9	Fc gamma R	04666	79	42	2.38e-05	-10.84	0.250000	7.75e-05	9.56e-04
10	MAPK signa	04010	246	104	6.34e-05	11.27	0.168000	1.33e-04	1.47e-03
11	Small cell	05222	75	35	2.45e-03	25.47	0.006000	1.78e-04	1.80e-03
12	Regulation	04810	192	84	7.58e-05	8.25	0.522000	4.41e-04	4.08e-03
13	Wnt signal	04310	138	60	9.08e-04	-6.67	0.339000	2.80e-03	2.39e-02
14	Bacterial	05100	62	32	4.36e-04	2.94	0.814000	3.17e-03	2.52e-02
15	Staphyloco	05150	43	20	2.03e-02	14.34	0.033000	5.56e-03	4.11e-02
16	Renal cell	05211	64	29	9.15e-03	-8.26	0.116000	8.33e-03	5.47e-02
17	Colorectal	05210	57	24	4.32e-02	8.45	0.025000	8.46e-03	5.47e-02
18	B cell rec	04662	70	32	5.43e-03	-10.21	0.210000	8.87e-03	5.47e-02
19	ErbB signa	04012	76	33	1.21e-02	-18.55	0.111000	1.02e-02	5.76e-02
20	Pathogenic	05130	48	21	3.69e-02	17.59	0.037000	1.04e-02	5.76e-02
	pGFWER	Status							
1	1.20e-09	Activated							
2	3.70e-08	Inhibited							
3	3.81e-08	Activated							
4	4.51e-08	Inhibited							
5	1.20e-04	Activated							
6	7.47e-04	Inhibited							
7	1.39e-03	Inhibited							
8	2.85e-03	Activated							
9	8.60e-03	Inhibited							
10	1.47e-02	Activated							
11	1.98e-02	Activated							
12	4.89e-02	Activated							
13	3.10e-01	Inhibited							
14	3.52e-01	Activated							
15	6.17e-01	Activated							

```

16 9.25e-01 Inhibited
17 9.39e-01 Activated
18 9.84e-01 Inhibited
19 1.00e+00 Inhibited
20 1.00e+00 Activated

```

If the `plots` argument is set to `TRUE` in the function call above, a plot like the one shown in Figure 1 is produced for each pathway on which there are differentially expressed genes. These plots are saved in a pdf file in the current directory.

An overall picture of the pathways significance according to both the over-representation evidence and perturbations based evidence can be obtained with the function `plotP` and shown in Figure 2. The Colorectal cancer pathway is shown in green.

In this plot, the horizontal axis represents the p-value (minus log of) corresponding to the probability of obtaining at least the observed number of genes (NDE) on the given pathway just by chance. The vertical axis represents the p-value (minus log of) corresponding to the probability of obtaining the observed total accumulation (tA) or more extreme on the given pathway just by chance. The computation of pPERT is described in Tarca et al. (2009). In Figure 2 each pathway is shown as a bullet point, and those significant at 5% (set by the `threshold` argument in `plotP`) after Bonferroni correction are shown in red.

The default method to combine pPERT and pNDE is Fisher's product method, as was described in Tarca et al. (2009).

Alternatively, the two types of evidence can be combined using a normal inversion method which gives smaller pG values when pPERT and pNDE are low simultaneously. This is in contrast with Fisher's method that may yield small pG values when only one of the two p-values is low. To use the normal inversion method, one can set the argument `combine="norminv"` when the `spia` function is called, or by recomputing pG values starting with a result data frame produced by `spia` function. This latter approach is illustrated below where a call is made to the function `combfunc`. SPIA algorithm is illustrated also using the Vessels dataset:

```

> data(Vessels)
> res <- spia(de = DE_Vessels, all = ALL_Vessels, organism = "hsa",
+   nB = 500, plots = FALSE, beta = NULL, verbose = FALSE)
> res$Name = substr(res$Name, 1, 10)
> res[1:15, -12]

```

	Name	ID	pSize	NDE	pNDE	tA	pPERT	pG	pGFdr	pGFWER
1	Axon guida	04360	129	12	2.24e-04	-6.6959	0.03600	0.000103	0.0105	0.0105
2	Olfactory	04740	378	1	1.00e+00	-4.0970	0.00002	0.000236	0.0106	0.0241
3	Staphyloco	05150	51	8	6.74e-05	2.4721	0.40000	0.000311	0.0106	0.0317
4	Focal adhe	04510	199	16	1.23e-04	-5.6255	0.38000	0.000514	0.0131	0.0524
5	Viral myoc	05416	70	8	6.35e-04	-1.6656	0.24800	0.001536	0.0254	0.1567
6	Rheumatoid	05323	89	10	1.59e-04	0.0000	1.00000	0.001547	0.0254	0.1578
7	Neuroactiv	04080	271	18	5.12e-04	-0.5104	0.41600	0.002015	0.0254	0.2055
8	Intestinal	04672	46	7	2.34e-04	0.0000	1.00000	0.002193	0.0254	0.2237
9	Regulation	04810	210	14	2.00e-03	7.2337	0.12000	0.002241	0.0254	0.2286

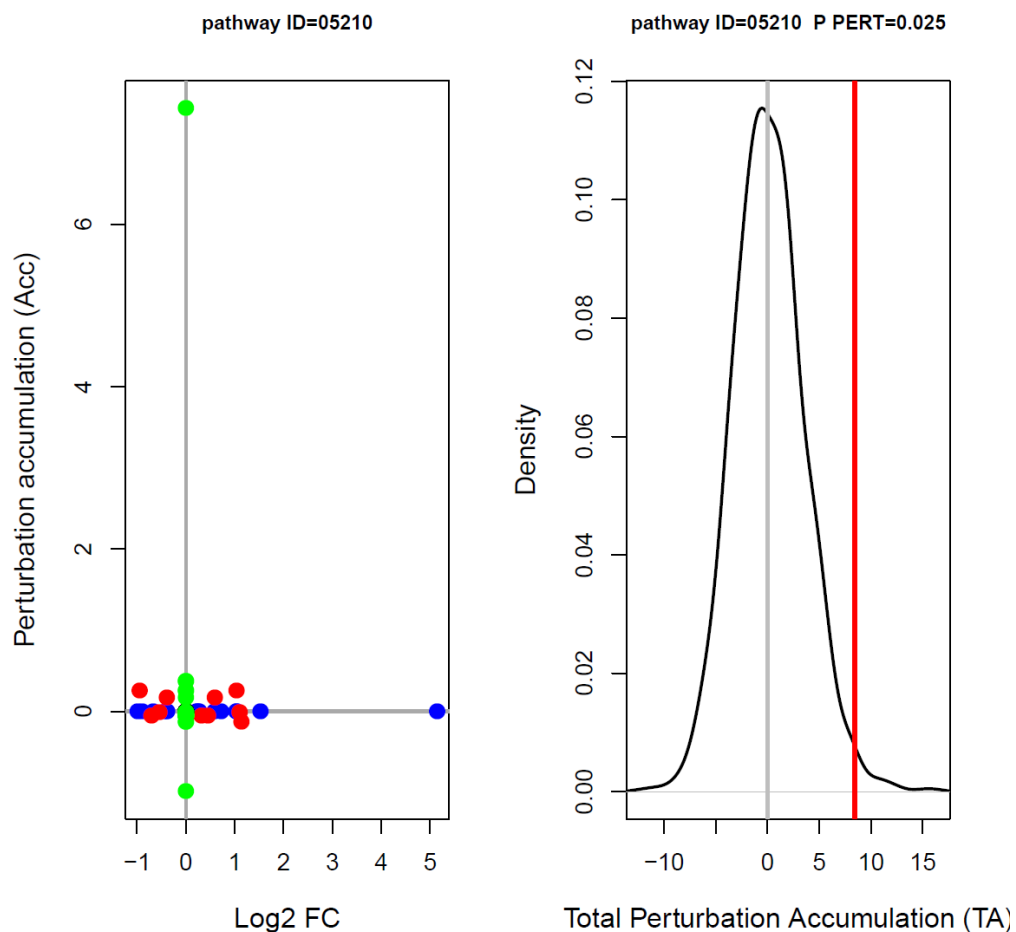


Figure 1: Perturbations plot for colorectal cancer pathway (KEGG ID hsa:05210) using the `colorectal_cancer` dataset. The perturbation of all genes in the pathway are shown as a function of their initial log2 fold changes (left panel). Non DE genes are assigned 0 log2 fold-change. The null distribution of the net accumulated perturbations is also given (right panel). The observed net accumulation tA with the real data is shown as a red vertical line.

```

> plotP(res, threshold = 0.05)
> points(I(-log(pPERT)) ~ I(-log(pNDE)), data = res[res$ID == "05210",
+           ], col = "green", pch = 19, cex = 1.5)

```

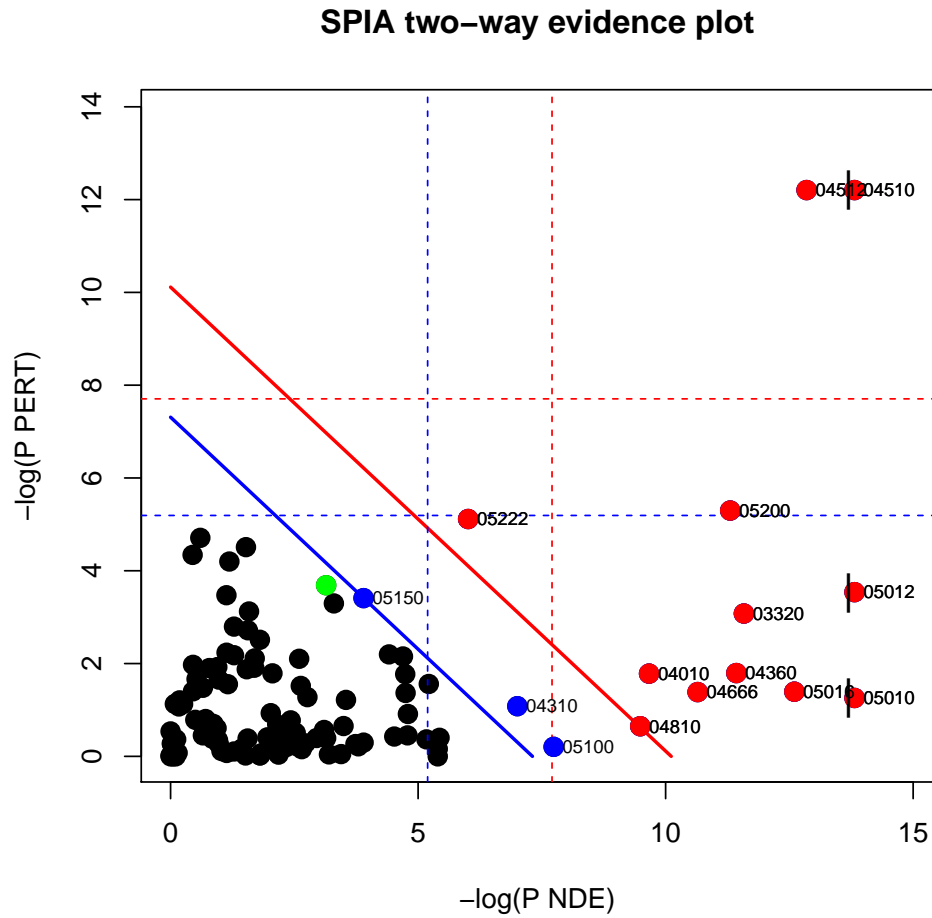


Figure 2: SPIA evidence plot for the colorectal cancer dataset. Each pathway is represented by one dot. The pathways at the right of the red oblique line are significant after Bonferroni correction of the global p-values, pG, obtained by combining the pPERT and pNDE using Fisher's method. The pathways at the right of the blue oblique line are significant after a FDR correction of the global p-values, pG.

```

> res$pG = combfunc(res$pNDE, res$pPERT, combine = "norminv")
> res$pGFdr = p.adjust(res$pG, "fdr")
> res$pGFWER = p.adjust(res$pG, "bonferroni")
> plotP(res, threshold = 0.05)
> points(I(-log(pPERT)) ~ I(-log(pNDE)), data = res[res$ID == "05210",
+      ], col = "green", pch = 19, cex = 1.5)

```

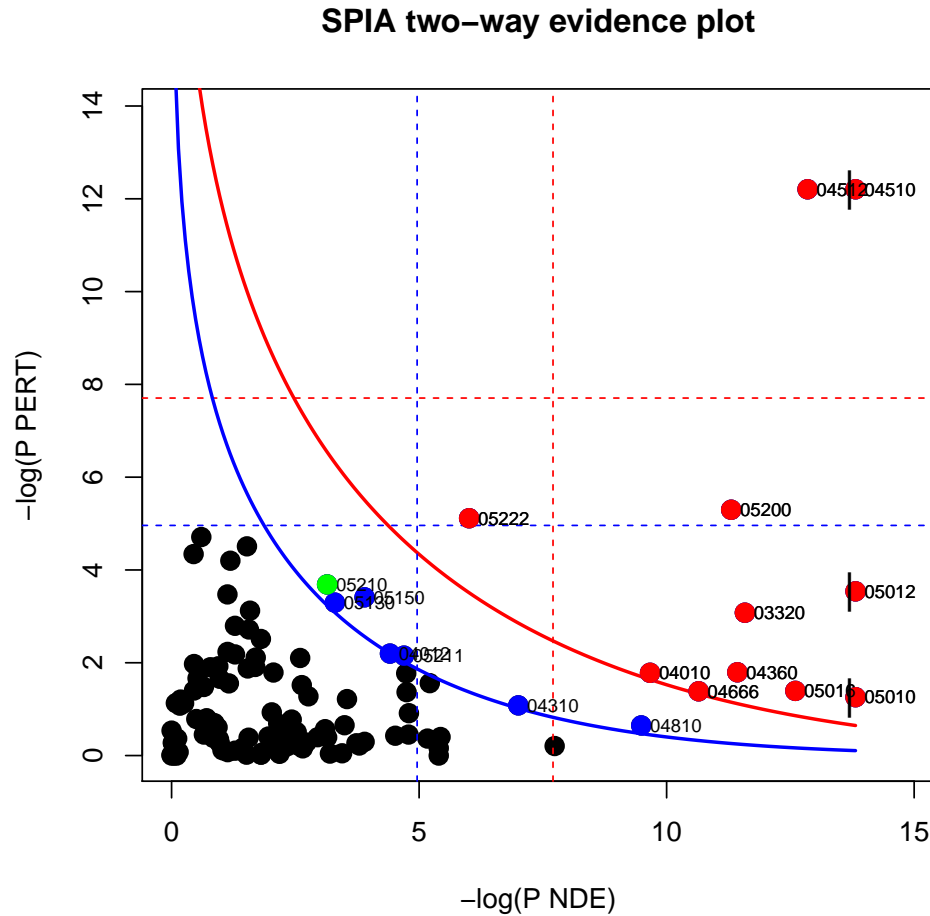


Figure 3: SPIA evidence plot for the colorectal cancer dataset. Each pathway is represented by one dot. The pathways at the right of the red curve are significant after Bonferroni correction of the global p-values, pG , obtained by combining the $pPERT$ and $pNDE$ using the normal inversion method. The pathways at the right of the blue curve line are significant after a FDR correction of the global p-values, pG .

10	Antigen pr	04612	76	7	4.74e-03	1.4922	0.09200	0.003809	0.0389	0.3886
11	Leishmania	05140	68	8	5.22e-04	-0.0638	0.98400	0.004400	0.0408	0.4488
12	Graft-vers	05332	41	6	8.13e-04	0.0000	1.00000	0.006597	0.0527	0.6729
13	Complement	04610	67	7	2.33e-03	4.7219	0.36400	0.006834	0.0527	0.6971
14	Notch sign	04330	46	4	3.66e-02	7.4277	0.02800	0.008079	0.0527	0.8241
15	Asthma	05310	29	5	1.04e-03	0.0000	1.00000	0.008167	0.0527	0.8330

Status

- 1 Inhibited
- 2 Inhibited
- 3 Activated
- 4 Inhibited
- 5 Inhibited
- 6 Inhibited
- 7 Inhibited
- 8 Inhibited
- 9 Activated
- 10 Activated
- 11 Inhibited
- 12 Inhibited
- 13 Activated
- 14 Activated
- 15 Inhibited

The pathway image as provided by KEGG having the differentially expressed genes highlighted in red can be obtained by pasting in a web browser the links available in the KEGGLINK column of the data frame produced by the function spia. For example,

```
> res[, "KEGGLINK"][20]
```

```
[1] "http://www.genome.jp/dbget-bin/show_pathway?hsa04110+983+7533+9232+595"
```

is the link that would display the image of the 20th pathway in the res dataframe above.

Note that the results for these datasets may differ from the ones described in Tarca et al. (2009) since a) the pathways database used herein was updated and b) the default beta values were changed.

The directed adjacency matrices of the graphs describing the different types of relations between genes/proteins (such as activation or repression) used by SPIA are available in the `extdata/hsaSPIA.RData` file for the homo sapiens organism. The types of relations considered by SPIA and the default weight (beta coefficient) given to them are:

```
> rel <- c("activation", "compound", "binding/association", "expression",
+ "inhibition", "activation_phosphorylation", "phosphorylation",
+ "indirect", "inhibition_phosphorylation", "dephosphorylation_inhibition",
+ "dissociation", "dephosphorylation", "activation_dephosphorylation",
+ "state", "activation_indirect", "inhibition_ubiquination",
+ "ubiquination", "expression_indirect", "indirect_inhibition",
+ "repression", "binding/association_phosphorylation", "dissociation_phosphorylation",
```



```

+     "indirect_phosphorylation")
> beta = c(1, 0, 0, 1, -1, 1, 0, 0, -1, -1, 0, 0, 1, 0, 1, -1,
+     0, 1, -1, -1, 0, 0, 0)
> names(beta) <- rel
> cbind(beta)

```

	beta
activation	1
compound	0
binding/association	0
expression	1
inhibition	-1
activation_phosphorylation	1
phosphorylation	0
indirect	0
inhibition_phosphorylation	-1
dephosphorylation_inhibition	-1
dissociation	0
dephosphorylation	0
activation_dephosphorylation	1
state	0
activation_indirect	1
inhibition_ubiquination	-1
ubiquination	0
expression_indirect	1
indirect_inhibition	-1
repression	-1
binding/association_phosphorylation	0
dissociation_phosphorylation	0
indirect_phosphorylation	0

A 0 value for a given relation type results in discarding those type of relations from the analysis for all pathways. The default values of **beta** can be changed by the user at any time by setting the **beta** argument of the **spia** function call.

Other organisms' KEGG pathway data can be downloaded from <http://bioinformaticsprb.med.wayne.edu/SPIA> as a "[org]SPIA.RData" file and copied into the **extdata** directory of the SPIA package, and therefore make it available to the function **spia**.

The user has the ability to generate his own gene/protein relation data and put it in a list format as the one shown in the **hsaSPIA.RData** file. In this file, each pathway data is included in a list:

```

> load(file = paste(system.file("extdata/hsaSPIA.RData", package = "SPIA")))
> names(path.info[["05210"]])

[1] "activation"           "compound"
[3] "binding/association"  "expression"
[5] "inhibition"           "activation_phosphorylation"
[7] "phosphorylation"      "indirect"

```

```

[9] "inhibition_phosphorylation"      "dephosphorylation_inhibition"
[11] "dissociation"                    "dephosphorylation"
[13] "activation_dephosphorylation"    "state"
[15] "activation_indirect"             "inhibition_ubiquination"
[17] "ubiquination"                   "expression_indirect"
[19] "indirect_inhibition"             "repression"
[21] "binding/association_phosphorylation" "dissociation_phosphorylation"
[23] "indirect_phosphorylation"        "nodes"
[25] "title"                           "NumberOfReactions"

```

```
> path.info[["05210"]][["activation"]][25:35, 30:40]
```

	5602	8312	8313	5900	387	5879	5880	5881	332	4609	595
369	0	0	0	0	0	0	0	0	0	0	0
5894	0	0	0	0	0	0	0	0	0	0	0
673	0	0	0	0	0	0	0	0	0	0	0
5599	0	0	0	0	1	1	1	1	0	0	0
5601	0	0	0	0	1	1	1	1	0	0	0
5602	0	0	0	0	1	1	1	1	0	0	0
8312	0	0	0	0	0	0	0	0	0	0	0
8313	0	0	0	0	0	0	0	0	0	0	0
5900	0	0	0	0	0	0	0	0	0	0	0
387	0	0	0	1	0	0	0	0	0	0	0
5879	0	0	0	1	0	0	0	0	0	0	0

In the matrix above, only 0 and 1 values are allowed. 1 means the gene/protein given by the column has a relation of type "activation" with the gene/protein given by the row of the matrix.

Using other R packages such as **graph** and **Rgraphviz** one can visualize the richness of gene/protein relations of each type in each pathway. Firstly we load the required packages and create a function that can be used to plot as a graph each type of relation of any pathway, as used by SPIA.

```

> library(graph)
> library(Rgraphviz)
> plotG <- function(B) {
+   nnms <- NULL
+   colls <- NULL
+   mynodes <- colnames(B)
+   L <- list()
+   n <- dim(B)[1]
+   for (i in 1:n) {
+     L[i] <- list(edges = rownames(B)[abs(B[, i]) > 0])
+     if (sum(B[, i] != 0) > 0) {
+       nnms <- c(nnms, paste(colnames(B)[i], rownames(B)[B[,
+         i] != 0], sep = "~"))
+     }
+   }
+ }

```

```

+   names(L) <- rownames(B)
+   g <- new("graphNEL", nodes = mynodes, edgeL = L, edgemode = "directed")
+   plot(g)
+ }

```

We plot then the "activation" relations in the ErbB signaling pathway, based on the `hsaSPIA` data.

```

> plotG(path.info[["04012"]][["activation"]])

```

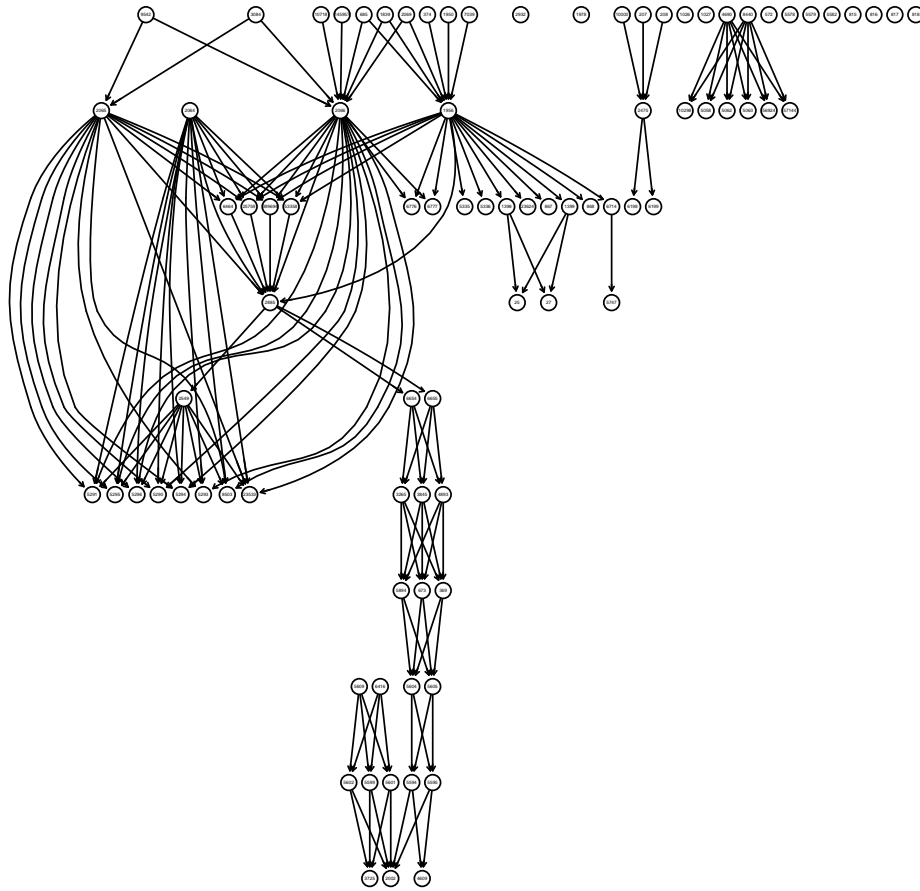


Figure 4: Display of the "activation" relations in the ErbB signaling pathway, based on the `hsaSPIA` data.

For more details on how to use the main function in this package use `"?spia"`.

3 Changes in SPIA 2.0 vs 1.9

The current version (2.0) contains the following changes compared to the previous version (1.9):

- 1) The computation of $pPERT$ was improved by replacing to NA the $pPERT$ of a pathway whenever the observed tA value was 0.0 and all the empirical distribution of tA values under the null hypothesis is made only of 0.0 values. In those rare cases, there is nothing to be learned from the tA values therefore instead of assigning $pPERT$ to 1.0, as in the past version, now we replace it with NA. Therefore the global pG value will be equal to the over-representation p-value ($pNDE$).
- 2) The function `getP2` used in the `plotP` is improved by using an analytical method instead of a numerical method to compute the probability $P1$ and $P2$ (with $P1 = P2$) such that the combined probability pG is equal to the desired threshold passed as argument to the function `pPlot`. This function is only used in creating the SPIA two way evidence plot.
- 3) The perturbation factor coming from a given upstream gene A used to be divided by the number of all downstream genes of A for each type of relation individually. In the current version, the division is made to the number of all downstream genes of gene A including all types of relations considered (i.e. having non zero β).
- 4) An additional meta-analysis method, call normal inversion was added for combining $pPERT$ and $pNDE$ into a global probability pG , in addition to the initial method (Fisher's product).
- 5) The date stamp for KEGG's pathway data is 03/21/2011.

References

- S. Draghici, P. Khatri, A. Tarca, K. Amin, A. Done, C. Voichita, C. Georgescu, and R. Romero. A systems biology approach for pathway level analysis. *Genome Research*, 17, 2007.
- Y. Hong, K. S. Ho, K. W. Eu, and P. Y. Cheah. A susceptibility gene set for early onset colorectal cancer that integrates diverse signaling pathways: implication for tumorigenesis. *Clin Cancer Res*, 13(4):1107–14, 2007.
- P. Khatri, S. Draghici, A. L. Tarca, S. S. Hassan, and R. Romero. A system biology approach for the steady-state analysis of gene signaling networks. In *12th Iberoamerican Congress on Pattern Recognition*, Valparaiso, Chile, November 13-16 2007.
- A. L. Tarca, S. Draghici, P. Khatri, S. Hassan, P. Mital, J. Kim, C. Kim, J. P. Kusanovic, and R. Romero. A signaling pathway impact analysis for microarray experiments. *Bioinformatics*, 25:75–82, 2009.