

Centrality-based Pathway Enrichment

Zuguang Gu

June 29, 2012

1 Introduction

Gene set enrichment analysis is broadly used in microarray data analysis [8, 5]. It aims to find which biological functions are affected by a group of related genes behind the massive information. The most used methodology is finding these significant gene set from a 2×2 contingency table, usually by Fisher's exact test or chi-square test. This kind of analysis is known as Over-represented Analysis (ORA). It takes a list of differential expressed gene, and returns significant gene sets that the differential genes are enriched in. A lot of methods have been developed under the framework of ORA such as DAVID [6] (<http://david.abcc.ncifcrf.gov/>) and `G0stats` package [3]. The second methodology to find significant pathways is to use whole expression matrix, named Gene-set Analysis (GSA). GSA methods are implemented via either a univariate or a multivariate procedure [1]. In univariate analysis, gene level statistics are initially calculated from fold changes or statistical tests (e.g., t -test). These statistics are then combined into a pathway level statistic by summation or averaging. GSEA [11] is a widely used univariate tool that utilizes a weighted Kolmogorov-Smirnov test to measure the degree of differential expression of a gene set by calculating a running sum from the top of a ranked gene list. Multivariate analysis considers the correlations between genes in the pathway and calculates the pathway level statistic directly from the expression value matrix using Hotelling's T^2 test [10] or MANOVA models [7].

For a specific form of gene sets, biological pathways are collections of correlated genes/proteins, RNAs and compounds that work together to regulate specific biological processes. Instead of just being a list of genes, a pathway contains the most important information that is how the member genes interact with each other. Thus network structure information is necessary for the interpretation of the importance of the pathways.

In this package, the original pathway enrichment method (ORA and GSA) is extended by introducing network centralities as the weight of nodes which have been mapped from differentially expressed genes in pathways [4]. There are two advantages compared to former methods. First, for the diversity of genes' characters and the difficulties of covering the importance of genes from all aspects, we do not design a fixed measurement for each gene but set it as an optional parameter in the model. Researchers can select from candidate choices where different measurement reflects different aspect of the importance of genes. In our model, network centralities are used to measure the importance of genes in pathways. Different centrality measurements assign the importance to nodes

from different aspects. For example, degree centrality measures the amount of neighbours that a node directly connects to, and betweenness centrality measures how many information streams must pass through a certain node. Generally speaking, nodes having large centrality values are central nodes in the network. It's observed that nodes represented as metabolites, proteins or genes with high centralities are essential to keep the steady state of biological networks. Moreover, different centrality measurements may relate to different biological functions. The selection of centralities for researchers depends on what kind of genes they think important. Second, we use nodes as the basic units of pathways instead of genes. We observe that nodes in the pathways include different types of molecules, such as single gene, complex and protein families. Assuming a complex or family contains ten differentially expressed member genes, in traditional ORA, these ten genes behave as the same position as other genes represented as single nodes, and thus they have effect of ten. It is not proper because these ten genes stay in a same node in the pathway and make functions with the effect of one node. Also, a same gene may locate in different complexes in a pathway and if taking the gene with effect of one, it would greatly decrease the importance of the gene. Therefore a mapping procedure from genes to pathway nodes is applied in our model. What's more, the nodes in pathways also include non-gene nodes such as microRNAs and compounds. These nodes also contribute to the topology of the pathway. So, when analyzing pathways, all types of nodes are retained.

2 Pathway Catalogue

Pathways are collected from public databases, such as PID, KEGG, BioCarta etc. In **CePa** package, four catalogues (PID, KEGG, BioCarta and Reactome) from PID database have been integrated. The pathway data are parsed from XML format file provided by the PID FTP site. The Perl code for parsing can be obtained from the author's website (<http://mcube.nju.edu.cn/jwang/lab/soft/cepa/>). The pathway data is stored in `PID.db`.

```
> library(CePa)
> data(PID.db)
> names(PID.db)
```

```
[1] "NCI"          "BioCarta" "KEGG"       "Reactome"
```

Each pathway catalogue has been stored as a `pathway.catalogue` class object. The `print.pathway.catalogue` function simply prints the number of pathways in the catalogue. The `plot.pathway.catalogue` function visualizes general information of the catalogue (figure 1). It plot: A) Distribution of the number of member genes in each node; B) Distribution of the number of nodes in which a single gene resides; C) Relationship between node count and gene count in biological pathways.

```
> class(PID.db$NCI)
```

```
[1] "pathway.catalogue"
```

```
> PID.db$NCI
```

The catalogue contains 225 pathways.

```
> plot(PID.db$NCI)
```

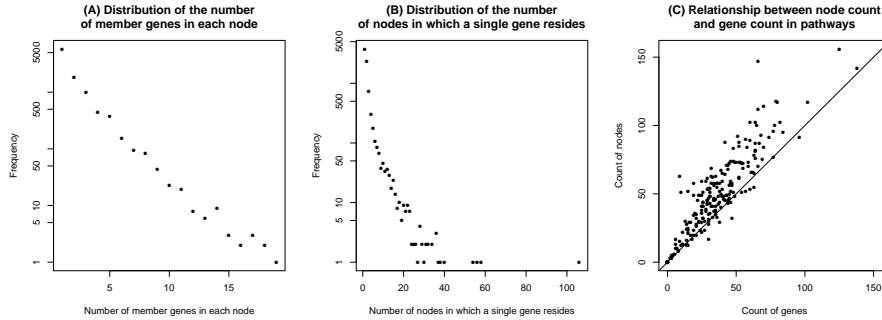


Figure 1: Meta analysis of pathway catalogue

The pathway catalogue data contains a list of pathways and each pathway contains a list of interactions. There are several parts in the pathway data where three of them is must: the pathway list, the interaction list and the mapping list. The corresponding list name are `pathList`, `interactionList` and `mapping`.

```
> names(PID.db$NCI)
```

```
[1] "pathList"          "interactionList" "mapping"          "node.name"
[5] "node.type"         "created"
```

The `pathList` is a list in which each item is a list of interaction IDs

```
> head(PID.db$NCI$pathList[[1]])
```

```
[1] "209941" "209933" "209957" "209927" "209931" "209958"
```

The `interactionList` is a three-column matrix in which the first column is the interaction ID, the second column is the input node ID and the third column is the output node ID.

```
> head(PID.db$NCI$interactionList)
```

```
interaction.id input output
1          201405 200665 205628
2          201405 200666 205628
3          204164 208481 208484
4          204164 202538 208484
5          206327 200592 200709
6          206327 210859 200709
```

The `mapping` is the two-column matrix in which the first column is the node ID and the second column is the gene ID.

```
> head(PID.db$NCI$mapping)
```

	node.id	symbol
1	201978	HGS
2	202230	ARHGAP6
3	201405	XIAP
4	201647	CRY2
5	202024	GZMA
6	201386	HFE2

The pathway catalogue can also be self-defined by `set.pathway.catalogue` function. The function returns a `pathway.catalogue` class object. E.g. we only need the first ten pathways in NCI catalogue.

```
> new.catalogue = set.pathway.catalogue(pathList = PID.db$NCI$pathList[1:10],
+                                       interactionList = PID.db$NCI$interactionList,
+                                       mapping = PID.db$NCI$mapping)
```

In the following examples, we will use NCI catalogue as the default pathway catalogue.

3 ORA Extension

The pathway score is defined as the summation of the weights of differentially affected nodes in the pathway:

$$s = \sum_{i=1}^n w_i d_i \quad (1)$$

where s is the score of the pathway, w_i is the weight of the i^{th} node and reflects the importance of the node, n is the number of nodes in the pathway, and d_i identifies whether the i^{th} node is differentially affected ($= 1$) or not ($= 0$).

The **CePa** package needs a differentially expressed gene list and a background gene list. The differential gene list can be obtained through variety of methods such as t -test, SAM [12] and limma [9]. The background gene list is the complete category of genes that exist on a certain microarray platform or from the whole genome. The **CePa** package contains an example gene list and a background gene list. The gene list is obtained from a microarray study by t -test [2].

```
> data(gene.list)
> names(gene.list)

[1] "bk" "dif"
```

In order to find significant pathways under several centrality measurements, we use `cepa.all` function. In the function, `dif` refers to the differential gene list, `bk` refers to the background gene list and the `pc` refers to the pathway catalogue.

```
> res = cepa.all(dif = gene.list$dif, bk = gene.list$bk,
+               pc = PID.db$NCI)
```

```

Calculate pathway scores...
1/205, hif1_tfp pathway...
- equal.weight: 0.7212787
- in.degree: 0.7972028
- out.degree: 0.8321678
- betweenness: 0.8031968
- in.reach: 0.6603397
- out.reach: 0.8141858
...

```

The differential gene list and the background gene list should be indicated with the same identifiers (e.g. gene symbol or refseq ID). All genes in the differential gene list should exist in the background gene list. In this example, gene list must be formatted as gene symbol. If background gene list is not specified, the function use whole human genome genes as default.

By default, `cepa.all` use `equal.weight`, `in.degree`, `out.degree`, `betweenness`, `in.reach` and `out.reach` centralities as pathway nodes' weight. More centrality measurements can be used by setting it as a function (such as closeness, cluster coefficient).

In order to generate the null distribution of the pathway score, novel differential gene list is sampled from the background gene list. P-values are calculated from 1000 simulations by default.

`res` is a `cepa.all` class object. To see the general information of this object:

```

> res

number of pathways: 205

Significant pathways (p.value <= 0.01):
      Number
equal.weight      18
in.degree         20
out.degree        16
betweenness       15
in.reach          21
out.reach         20

```

It will print the number of significant pathways under different centralities. For ORA extension, `cepa.all` in fact calls `cepa.ora.all` function. So the following code is same as the former code.

```

> res = cepa.ora.all(dif = gene.list$dif, bk = gene.list$bk,
+                   pc = PID.db$NCI)

```

The p-values or adjusted p-values of all pathways under different centralities can be compared through the heatmap of p-values (Figure 2). Users can select methods to adjust raw p-values.

```

> plot(res, adj.method = "BH", cutoff = 0.05)

```

By default, `plot` generates the heatmap containing all pathways. If only significant pathways are of interest, the `only.sig` argument can be set to `TRUE`. (Figure 3).

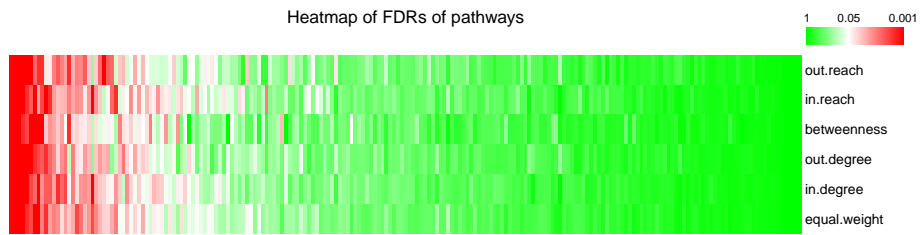


Figure 2: Heatmap of p-values of all pathways

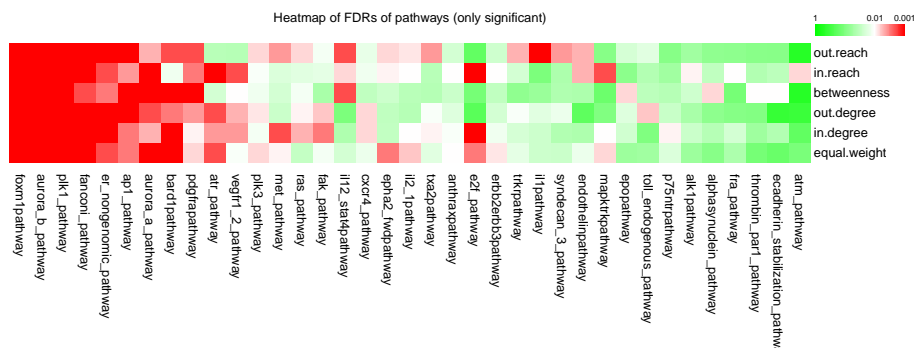


Figure 3: Heatmap of p-values of significant pathways

```
> plot(res, adj.method = "BH", only.sig = TRUE, cutoff = 0.01)
```

The numeric values of p-values can be obtained via `p.table`. The function just returns the raw p-values.

```
> pt = p.table(res)
> head(pt)
```

	equal.weight	in.degree	out.degree	betweenness
hif1_tfp_pathway	0.721278721	0.797202797	0.832167832	0.803196803
s1p_s1p5_pathway	0.144855145	0.199800200	0.160839161	0.055944056
wnt_signaling_pathway	0.869130869	0.864135864	0.907092907	0.896103896
ap1_pathway	0.002997003	0.002997003	0.000999001	0.000999001
lpa4_pathway	0.811188811	0.830169830	0.815184815	0.714285714
avb3_opn_pathway	0.796203796	0.594405594	0.649350649	0.310689311

	in.reach	out.reach
hif1_tfp_pathway	0.660339660	0.814185814
s1p_s1p5_pathway	0.069930070	0.145854146
wnt_signaling_pathway	0.809190809	0.906093906
ap1_pathway	0.003996004	0.000999001
lpa4_pathway	0.765234765	0.751248751
avb3_opn_pathway	0.729270729	0.589410589

We can get the result for single pathway under specific centrality from the `cepa.all` object by identifying the index for the pathway and the index for the centrality.

```
> g = get.cepa(res, id = "mapktrkpathway", cen = "in.degree")
> g

procedure: ora
weight: in.degree
p-value: 0.010
```

`g` is a `cepa` class object. It stores information of the evaluation of a single pathway under a single centrality. The distribution of the pathway score and the network graph can be generated by `plot` function on the `cepa` object by specifying `type` argument (figure 4 and figure 5).

```
> plot(g, type = "graph")
```

```
> plot(g, type = "null")
```

By default, `type` is set to `graph`, and the node labels is combined from member genes. The exact name for each node can be set by `node.name` argument. Also, more detailed categories of the nodes can be set by `node.type` argument (Figure 6).

```
> plot(g, node.name = PID.db$NCI$node.name,
+       node.type = PID.db$NCI$node.type)
```

For simplicity, the plotting for the `cepa` object can be directly applied on the `cepa.all` object by specifying the index of the pathway and the index of the centrality (Figure 6).

Graph view of the pathway

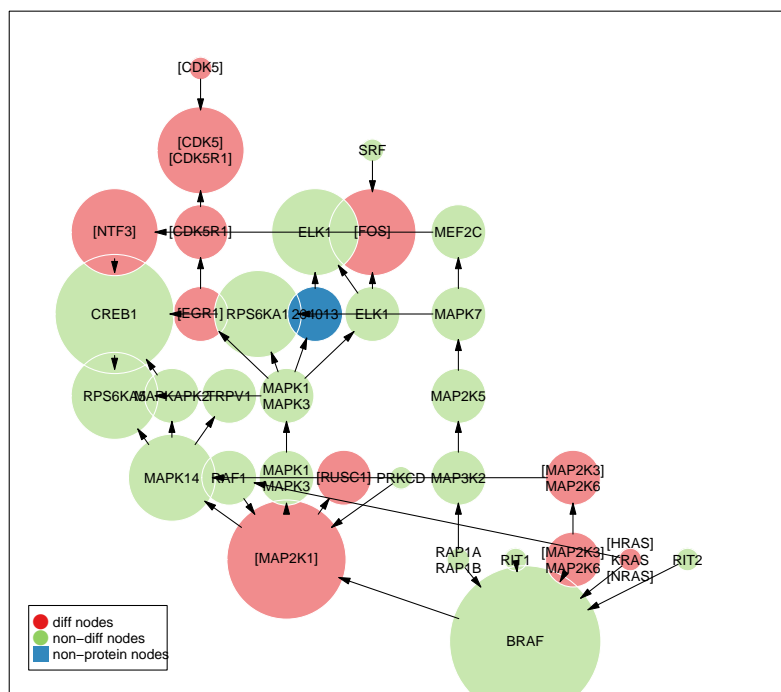


Figure 4: Network visualization of a pathway

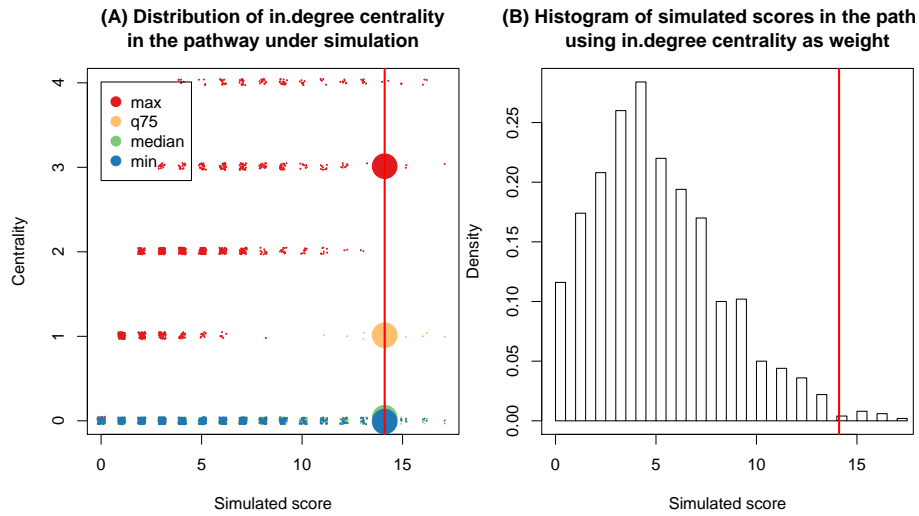


Figure 5: Null distribution of pathway score

```
> plot(res, id = "mapktrkpathway", cen = "in.degree")
> plot(res, id = "mapktrkpathway", cen = "in.degree", type = "null")
> plot(res, id = "mapktrkpathway", cen = "in.degree",
+       node.name = PID.db$NCI$node.name,
+       node.type = PID.db$NCI$node.type)
```

If used use `plot` to draw network graphs, the function would return an `igraph` object. So if users are not satisfy with the default graph, they can visulize by their own methods.

```
> obj = plot(res, id = "mapktrkpathway", cen = "in.degree")
> class(obj)
[1] "igraph"
```

The `igraph` package provides a `write.graph` function to output graph into several formats. As I have tried, with `graphml` format, Cytoscape Web (<http://http://cytoscapeweb.cytoscape.org/>) can make a more beautiful visualization of the network.

```
> write.graph(obj, file = "example-network.xml", format = "graphml")
> write.graph(obj, file = "example-network.gml", format = "gml")
```

Instead of analysis a list of pathways, users can also be focused on a single pathway under a single centrality by identifying the id of the pathway in the catalogue.

```
> res.pathway = cepa(dif = gene.list$dif, bk = gene.list$bk,
+                   pc = PID.db$NCI, id = 2)
```

Similarly, `cepa` function here directly calls `cepa.ora`.

Graph view of the pathway

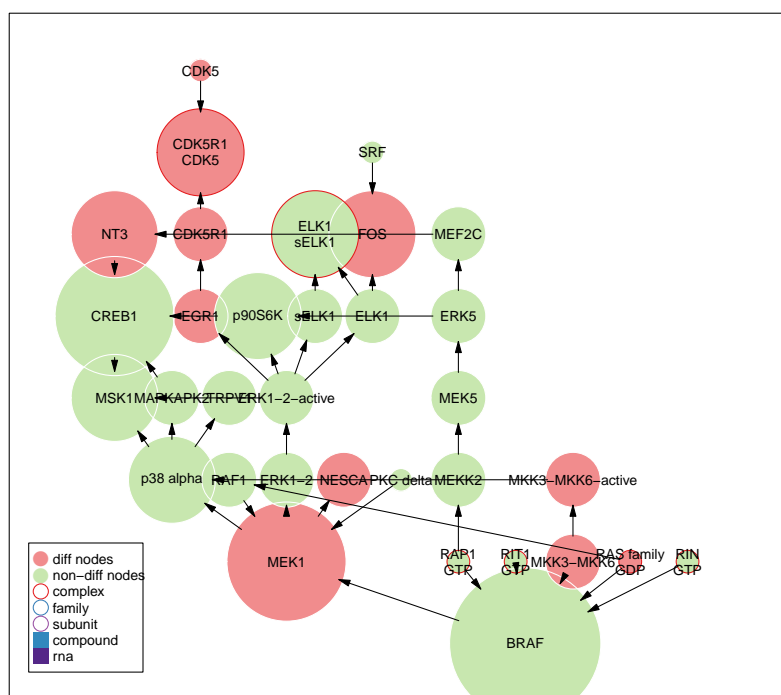


Figure 6: Network visualization of a pathway, with node name and node type specified

4 GSA extension

In the traditional univariate GSA procedure, the score s of the pathway is defined as:

$$s = f(\mathbf{g}) \quad (2)$$

where f transforms the gene-level statistic to a pathway-level statistic (e.g. by summation, averaging) and \mathbf{g} is the gene-level statistic vector which typically comprises t -values. In ORA, \mathbf{g} is a binary variant and $f(\mathbf{g})$ is summation. In our model to extend GSA, gene-level statistic is first transformed to node-level statistic. We define the vector of the node-level statistics as \mathbf{d} . When nodes in pathways comprise multiple genes, the node-level statistic can be considered as the largest principle component of the corresponding member genes. Using centrality as the weight, the score is defined as

$$s = f(\mathbf{w}\mathbf{d}) \quad (3)$$

where \mathbf{w} is the weight vector and the transformation function f acts upon the product of \mathbf{w} and \mathbf{d} . Equation 3 incorporates centrality weight into the original node-level statistic. The null distribution of the pathway score could then be generated by permuting the gene expression matrix.

Since GSA procedure need a complete expression matrix, we first read the P53 microarray data set. The `P53_symbol.gct` and `P53.cls` can be downloaded from <http://mcube.nju.edu.cn/jwang/lab/soft/cepa/>.

```
> eset = read.gct("P53_symbol.gct")
> # some process of the names of genes
> rownames(eset) = gsub("\\s+.*$", "", rownames(eset))
> label = read.cls("P53.cls", treatment="MUT", control="WT")
```

Here, we also use `cepa.all` to do batch pathway analysis.

```
> res = cepa.all(mat = eset, label = label, pc = PID.db$NCI,
               glevel = "tvalue_sq", plevel = "mean")
```

Calculate gene level values.

Calculate pathway score...

1/205, hif1_tfpathway...

Calculate node level value and permute sample labels...

- equal.weight: 0.788

- in.degree: 0.653

- out.degree: 0.405

- betweenness: 0.303

- in.reach: 0.917

- out.reach: 0.499

...

Here, we use `mat` and `label` arguments instead of `dif` and `bk` arguments. In fact, when specifying `mat` and `label` arguments, `cepa.all` calls `cepa.univariate.all`.

In GSA procedure, first a node level statistic should be calculated. In `CePa` package, there are three methods to calculate node level statistics. User can choose from `tvalue`, `tvalue_abs` and `tvalue_sq`. `tvalue_abs` is chosen as

the default node level method because it can capture two directional regulations. After we get the node level statistics in the pathway, a pathway level transformation should be applied. User can choose from `max`, `min`, `median`, `sum`, `mean` and `rank`. `mean` is taken as default.

Print the general result of the analysis and plot figures (figure 7).

```
> res

number of pathways: 205

Significant pathways (p.value <= 0.01):
      Number
equal.weight      5
in.degree         5
out.degree        7
betweenness       5
in.reach          6
out.reach         5

> plot(res, only.sig = TRUE, adj.method = "BH", cutoff = 0.15)
```

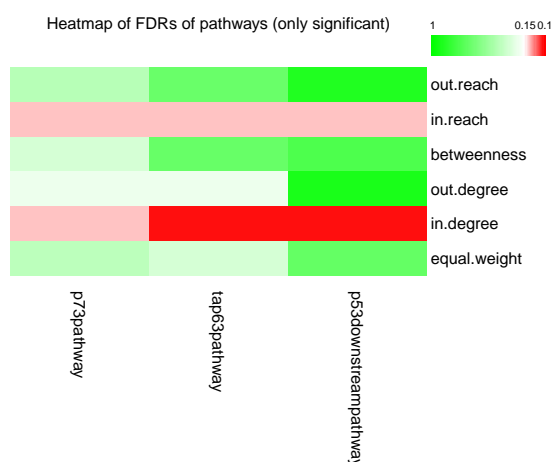


Figure 7: Heatmap of p-values of significant pathways

If we are instread in p73 pathway. First we extract this pathway under "in.degree" centrality from `res`.

```
> g = get.cepa(res, id = "p73pathway", cen="in.degree")
> g

procedure: gsa.univariate
weight: in.degree
p-value: 0.002
```

Graph view of the pathway

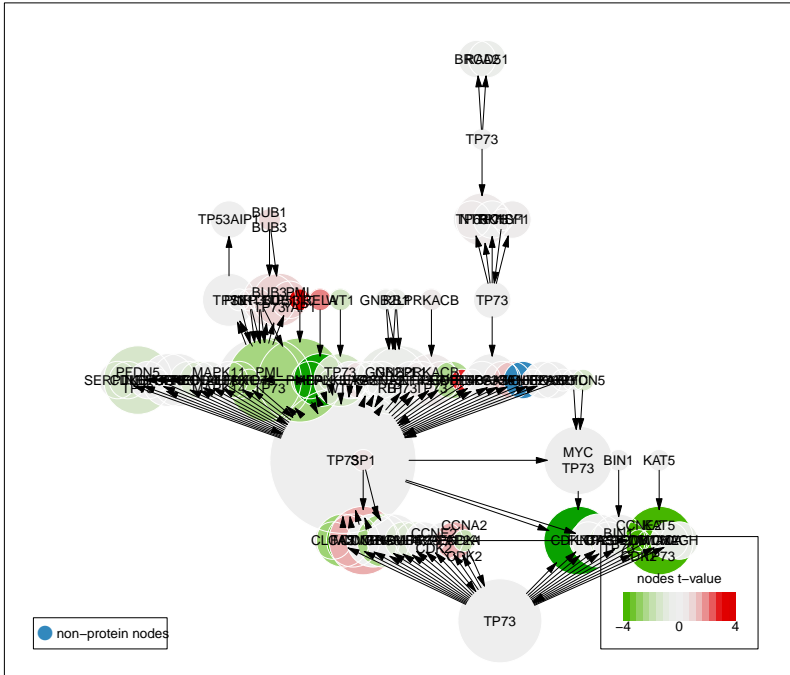


Figure 8: Network visulization of a pathway

```
> plot(g)
```

Figure 8 illustrates the graph of p73 pathway. Since the pathway is evaluated under GSA procedure, the color of each node is continues in which red refers to up-regulated, green refers to down-regulated and white refers to no-change.

5 The report function

One of the advantages of **CePa** package is that it can generate a detailed report in HTML format. The function **report** is used to generate report. The report will locate in the current working directory. By default it only generate figures of the significant pathways, but this can be changed by setting **only.sig** argument to **FALSE**.

```
> report(res)

generate images for ap1_pathway ...
generate images for epopathway ...
generate images for il12_stat4pathway ...
generate images for foxmlpathway ...
generate images for mapktrkpathway ...
generate images for aurora_a_pathway ...
...

> report(res, adj.method = "BH", cutoff = 0.15)
> report(res, sig.only = FALSE)
```

An example of the report can be found in figure 9.

6 Parallel computing

Since **CePa** evaluates pathways independently, the process can be realized through parallel computing. In R statistical environment, there are many packages focusing on parallel computing such as **snow**, **multicore**, etc. Here we demonstrate how to apply the parallel version of **CePa**, taking **multicore** for example.

```
> library(multicore)
> # identify how many cores you want to use in your computer
> ncores = 4
```

Since there are a list of pathways, we would link to divide them into several approximately equal groups, so we have a **divide** function (maybe you have a better function like this).

```
> divide = function(x, k) {
+   if(length(x) == 1 && is.numeric(x)) {
+     x = 1:x
+   }
+   if(length(x) < k) {
+     stop("o")
+   }
+ }
```

Analysis results by CePa
Method to adjust raw p-values: none
Cutoff: 0.01

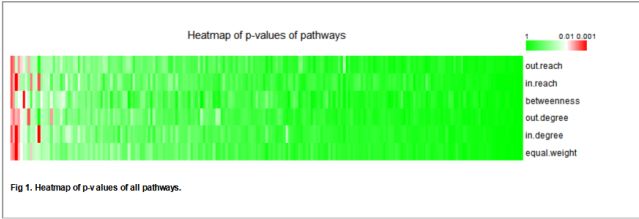


Fig 1. Heatmap of p-values of all pathways.

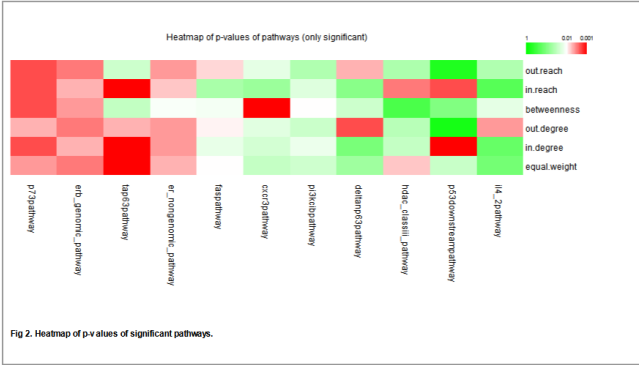


Fig 2. Heatmap of p-values of significant pathways.

Table 1. Complete list of p-values of all pathways under all centrality measurements.

Table of significance in b4 format

Significant pathways

Significant p-values

Pathway	equal.weight	in.degree	out.degree	betweenness	inreach	outreach
p73pathway	0.004	0.002	0.005	0.002	0.002	0.002
erb_genomic_pathway	0.003	0.005	0.003	0.004	0.005	0.003
tap63pathway	0.001	0.001	0.005	0.03	0.001	0.025
er_nongenomic_pathway	0.005	0.004	0.004	0.011	0.006	0.004
faspathway	0.01	0.015	0.009	0.012	0.047	0.007
cxc3pathway	0.029	0.022	0.017	0.001	0.0599	0.016
p3kic1bpathway	0.024	0.014	0.026	0.01	0.018	0.041
deltap63pathway	0.0559	0.1119	0.002	0.025	0.0839	0.005
hdac_classiii_pathway	0.006	0.029	0.036	0.2617	0.003	0.044
p53downstreampathway	0.027	0.001	0.6823	0.0959	0.002	0.5654
il4_2pathway	0.1219	0.1570	0.004	0.016	0.2068	0.041
atb1_a6b4_integrin_pathway	0.032	0.039	0.036	0.017	0.031	0.034
il27pathway	0.025	0.018	0.0569	0.0649	0.026	0.032
c1rc.adipathway	0.0549	0.033	0.03	0.011	0.0649	0.0599

Figure 9: An report of the CePa analysis

```

+   }
+   w = floor(length(x)/k)
+   q = length(x) - k*w
+   d = matrix(0, nrow=k, ncol=2)
+   n = 1
+   for(i in 1:k) {
+     d[i, 1] = n
+     d[i, 2] = n+w-1+ifelse(q>0, 1, 0)
+     n = d[i,2]+1
+     q = ifelse(q > 0, q-1, 0)
+   }
+   d[k,2] = length(x)
+   return(d)
+ }

```

In the `divide` function, the first argument is a vector, usually a index vector, and the second argument identify how many part you want to divide into. Also, the first argument can be a positive integer. For example, we want to divide `1:10` into two groups.

```
> divide(1:10, 2)
```

```

      [,1] [,2]
[1,]     1     5
[2,]     6    10

```

The function returns a matrix. Rows correspond to groups and columns correspond to the start index and the end index. If the vector can not be divided equally, the function would return an approximately division.

```
> divide(1:10, 3)
```

```

      [,1] [,2]
[1,]     1     4
[2,]     5     7
[3,]     8    10

```

Now we can divide the complete NCI pathway catalogue into several groups.

```

> NCI = PID.db$NCI
> d = divide(1:length(NCI$pathList), ncores)

```

Then we use `mclapply` which is something like a parallel version of `lapply` to do parallele computing.

```

> res = mclapply(1:ncores, function(i) {
+   pc = NCI
+   pc$pathList = pc$pathList[d[i, 1]:d[i, 2]]
+   cepa.all(dif = dif, bk = bk, pc = pc)},
+   mc.cores = ncores)
> res = mclapply(1:ncores, function(i) {
+   pc = set.pathway.catalogue(pathList = PID.db$NCI$pathList[d[i, 1]:d[i, 2]],
+     interactionList = PID.db$NCI$interactionList,

```



```

                                mapping = PID.db$NCI$mapping)
+       cepa.all(mat = eset, label = label, pc = pc)},
+       mc.cores = ncores)

```

In the `mclapply`, calculation in each core would return a `cepa.all` object. Thus, `res` is a list of `cepa.all` objects. We need some code to transform it into a single `cepa.all` object containing all pathways.

```

> obj = list()
> for(i in 1:length(res)) {
+   obj = c(obj, res[[i]])
+ }
> class(obj) = "cepa.all"

```

OK, now the `obj` is a `cepa.all` object just like the one generated from non-parallel CePa.

References

- [1] M. Ackermann and K. Strimmer. A general modular framework for gene set enrichment analysis. *BMC bioinformatics*, 10:47, Jan. 2009.
- [2] J. Burchard, C. Zhang, A. M. Liu, R. T. P. Poon, N. P. Y. Lee, K.-F. Wong, P. C. Sham, B. Y. Lam, M. D. Ferguson, G. Tokiwa, R. Smith, B. Leeson, R. Beard, J. R. Lamb, L. Lim, M. Mao, H. Dai, and J. M. Luk. microRNA-122 as a regulator of mitochondrial metabolic gene network in hepatocellular carcinoma. *Molecular systems biology*, 6:402, Aug. 2010.
- [3] S. Falcon and R. Gentleman. Using GOstats to test gene lists for GO term association. *Bioinformatics*, 23(2):257–8, 2007.
- [4] Z. Gu, J. Liu, K. Cao, J. Zhang, and J. Wang. Centrality-based pathway enrichment: a systematic approach for finding significant pathways dominated by key genes. *BMC Systems Biology*, 6(1):56, 2012.
- [5] D. W. Huang, B. T. Sherman, and R. A. Lempicki. Bioinformatics enrichment tools: paths toward the comprehensive functional analysis of large gene lists. *Nucleic acids research*, 37(1):1–13, Jan. 2009.
- [6] D. W. Huang, B. T. Sherman, and R. A. Lempicki. Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nature protocols*, 4(1):44–57, Jan. 2009.
- [7] M. Hummel, R. Meister, and U. Mansmann. GlobalANCOVA: exploration and assessment of gene group effects. *Bioinformatics*, 24(1):78–85, Jan. 2008.
- [8] P. Khatra and S. Drăghici. Ontological analysis of gene expression data: current tools, limitations, and open problems. *Bioinformatics*, 21(18):3587–95, Sept. 2005.
- [9] G. K. Smyth. Limma: linear models for microarray data. pages 397–420, 2005.

- [10] S. Song and M. Black. *pcot2: Principal Coordinates and Hotelling's T-Square method*. R package version 1.24.0.
- [11] A. Subramanian, P. Tamayo, V. K. Mootha, S. Mukherjee, B. L. Ebert, M. A. Gillette, A. Paulovich, S. L. Pomeroy, T. R. Golub, E. S. Lander, and J. P. Mesirov. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences of the United States of America*, 102(43):15545–50, Oct. 2005.
- [12] V. G. Tusher, R. Tibshirani, and G. Chu. Significance analysis of microarrays applied to the ionizing radiation response. *Proceedings of the National Academy of Sciences of the United States of America*, 98(9):5116–21, Apr. 2001.