

Using *clusterProfiler* to identify and compare functional profiles of gene lists

Guangchuang Yu
School of Biological Sciences
The University of Hong Kong, Hong Kong SAR, China
email: gcyu@connect.hku.hk

April 11, 2014

Contents

1	Introduction	1
2	Citation	2
3	Gene Ontology Classification	2
4	Enrichment Analysis	3
4.1	Hypergeometric model	3
4.2	GO enrichment analysis	3
4.3	KEGG pathway enrichment analysis	4
4.4	DO enrichment analysis	5
4.5	Reactome pathway enrichment analysis	5
4.6	Function call	5
4.7	Visualization	5
4.7.1	barplot	6
4.7.2	cnetplot	7
4.7.3	pathview from pathview package	9
5	Biological theme comparison	10
6	Session Information	12

1 Introduction

In recently years, high-throughput experimental techniques such as microarray, RNA-Seq and mass spectrometry can detect cellular moleculars at systems-level. These kinds of analyses generate huge quantities of data, which need to be

given a biological interpretation. A commonly used approach is via clustering in the gene dimension for grouping different genes based on their similarities [1].

To search for shared functions among genes, a common way is to incorporate the biological knowledge, such as Gene Ontology (GO) and Kyoto Encyclopedia of genes and Genomes (KEGG), for identifying predominant biological themes of a collection of genes.

After clustering analysis, researchers not only want to determine whether there is a common theme of a particular gene cluster, but also to compare the biological themes among gene clusters. The manual step to choose interesting clusters followed by enrichment analysis on each selected cluster is slow and tedious. To bridge this gap, we designed *clusterProfiler* [2], for comparing and visualizing functional profiles among gene clusters.

2 Citation

Please cite the following articles when using *clusterProfiler*.

G Yu, LG Wang, Y Han, QY He. *clusterProfiler*: an R package for comparing biological themes among gene clusters. *OMICS: A Journal of Integrative Biology*. 2012, 16(5), 284-287.

3 Gene Ontology Classification

In *clusterProfiler*, `groupGO` is designed for gene classification based on GO distribution at a specific level.

```
require(DOSE)
data(geneList)
gene <- names(geneList)[abs(geneList) > 2]
head(gene)

## [1] "4312" "8318" "10874" "55143" "55388" "991"

ggo <- groupGO(gene = gene, organism = "human", ont = "BP",
               level = 3, readable = TRUE)
head(summary(ggo))

##              ID              Description Count
## GO:0019953 GO:0019953      sexual reproduction      9
## GO:0019954 GO:0019954      asexual reproduction      0
## GO:0022414 GO:0022414      reproductive process     18
```

```
## GO:0032504 GO:0032504      multicellular organism reproduction      10
## GO:0032505 GO:0032505      reproduction of a single-celled organism      0
## GO:0048610 GO:0048610      cellular process involved in reproduction      8
##
## GO:0019953                                     ASPM/CDK1/TRIP13
## GO:0019954
## GO:0022414 ASPM/CDK1/TRIP13/ID01/CCNB1/CSN3/PTTG1/COL16A1/DACH1/CORIN/GAMT/BMP4/
## GO:0032504                                     ASPM/TRIP13/CCNB1/CS
## GO:0032505
## GO:0048610                                     CDC20/TOP2A/
```

4 Enrichment Analysis

4.1 Hypergeometric model

Enrichment analysis [3] is a widely used approach to identify biological themes. Here we implement hypergeometric model to assess whether the number of selected genes associated with disease is larger than expected.

To determine whether any terms annotate a specified list of genes at frequency greater than that would be expected by chance, *clusterProfiler* calculates a p-value using the hypergeometric distribution:

$$p = 1 - \sum_{i=0}^{k-1} \frac{\binom{M}{i} \binom{N-M}{n-i}}{\binom{N}{n}}$$

In this equation, N is the total number of genes in the background distribution, M is the number of genes within that distribution that are annotated (either directly or indirectly) to the node of interest, n is the size of the list of genes of interest and k is the number of genes within that list which are annotated to the node. The background distribution by default is all the genes that have annotation.

P-values were adjusted for multiple comparison, and q-values were also calculated for FDR control.

4.2 GO enrichment analysis

```
ego <- enrichGO(gene = gene, universe = names(geneList),
  organism = "human", ont = "CC", pvalueCutoff = 0.01,
  readable = TRUE)
head(summary(ego))
```

```
## ID Description GeneRatio
## GO:0005819 GO:0005819 spindle 24/195
```

```

## GO:0015630 GO:0015630 microtubule cytoskeleton 37/195
## GO:0005876 GO:0005876 spindle microtubule 10/195
## GO:0000793 GO:0000793 condensed chromosome 16/195
## GO:0000779 GO:0000779 condensed chromosome, centromeric region 12/195
## GO:0044430 GO:0044430 cytoskeletal part 42/195
##
## BgRatio pvalue p.adjust qvalue
## GO:0005819 203/11855 2.74e-14 2.54e-12 1.61e-12
## GO:0015630 681/11855 8.69e-11 4.04e-09 2.56e-09
## GO:0005876 39/11855 4.83e-10 1.50e-08 9.50e-09
## GO:0000793 139/11855 9.98e-10 2.32e-08 1.47e-08
## GO:0000779 71/11855 1.53e-09 2.85e-08 1.81e-08
## GO:0044430 968/11855 4.00e-09 6.20e-08 3.93e-08
##
## GO:0005819
## GO:0015630 KIF20A/TACC3/CENPE/CHEK1/KIF18B/SKA1/TP
## GO:0005876
## GO:0000793
## GO:0000779
## GO:0044430 KIF20A/TACC3/CENPE/CHEK1/KIF18B/SKA1/TPX2/PSD3/KIF4A/ASPM/AK5/BIRC5/I
##
## Count
## GO:0005819 24
## GO:0015630 37
## GO:0005876 10
## GO:0000793 16
## GO:0000779 12
## GO:0044430 42

```

4.3 KEGG pathway enrichment analysis

```

kk <- enrichKEGG(gene = gene, organism = "human", pvalueCutoff = 0.01,
  readable = TRUE)
head(summary(kk))

```

```

## ID Description GeneRatio BgRatio
## hsa04110 hsa04110 Cell cycle 11/74 128/5894
## hsa04114 hsa04114 Oocyte meiosis 10/74 114/5894
## hsa03320 hsa03320 PPAR signaling pathway 7/74 70/5894
## hsa04914 hsa04914 Progesterone-mediated oocyte maturation 6/74 87/5894
## hsa04062 hsa04062 Chemokine signaling pathway 8/74 189/5894
## hsa04060 hsa04060 Cytokine-cytokine receptor interaction 9/74 265/5894
##
## pvalue p.adjust qvalue
## hsa04110 4.31e-07 3.02e-06 4.54e-07
## hsa04114 1.25e-06 4.38e-06 6.59e-07
## hsa03320 2.35e-05 5.49e-05 8.25e-06
## hsa04914 7.21e-04 1.26e-03 1.90e-04

```

##	hsa04062	2.37e-03	3.32e-03	5.00e-04		
##	hsa04060	5.58e-03	6.51e-03	9.79e-04		
##					geneID	Count
##	hsa04110	CDC45/CDC20/CCNB2/CCNA2/CDK1/MAD2L1/TTK/CHEK1/CCNB1/MCM5/PTTG1				11
##	hsa04114	CDC20/CCNB2/CDK1/MAD2L1/CALML5/AURKA/CCNB1/PTTG1/ITPR1/PGR				10
##	hsa03320	MMP1/FADS2/ADIPOQ/PCK1/FABP4/HMGCS2/PLIN1				7
##	hsa04914	CCNB2/CCNA2/CDK1/MAD2L1/CCNB1/PGR				6
##	hsa04062	CXCL10/CXCL13/CXCL11/CXCL9/CCL18/CCL8/CXCL14/CX3CR1				8
##	hsa04060	CXCL10/CXCL13/CXCL11/CXCL9/CCL18/IL1R2/CCL8/CXCL14/CX3CR1				9

4.4 DO enrichment analysis

Disease Ontology (DO) enrichment analysis is implemented in *DOSE*, please refer to the package vignettes. The `enrichDO` function is very useful for identifying disease association of interesting genes.

4.5 Reactome pathway enrichment analysis

With the demise of KEGG (at least without subscription), the KEGG pathway data in Bioconductor will not update and we encourage user to analyze pathway using *ReactomePA* which use Reactome as a source of pathway data. The function call of `enrichPathway` in *ReactomePA* is consistent with `enrichKEGG`.

4.6 Function call

The function calls of `groupGO`, `enrichGO`, `enrichKEGG`, `enrichDO` and `enrichPathway` are consistent. The input parameters of *gene* is a vector of entrezgene (for human and mouse) or ORF (for yeast) IDs, and *organism* should be supported species (please refer to the manual of the specific function).

For GO analysis, *ont* must be assigned to one of "BP", "MF", and "CC" for biological process, molecular function and cellular component, respectively. In `groupGO`, the *level* specify the GO level for gene projection.

In enrichment analysis, the *pvalueCutoff* is to restrict the result based on their pvalues and the adjusted p values. *Q-values* were also calculated for controlling false discovery rate (FDR).

The *readable* is a logical parameter to indicate the input gene IDs will map to gene symbols or not.

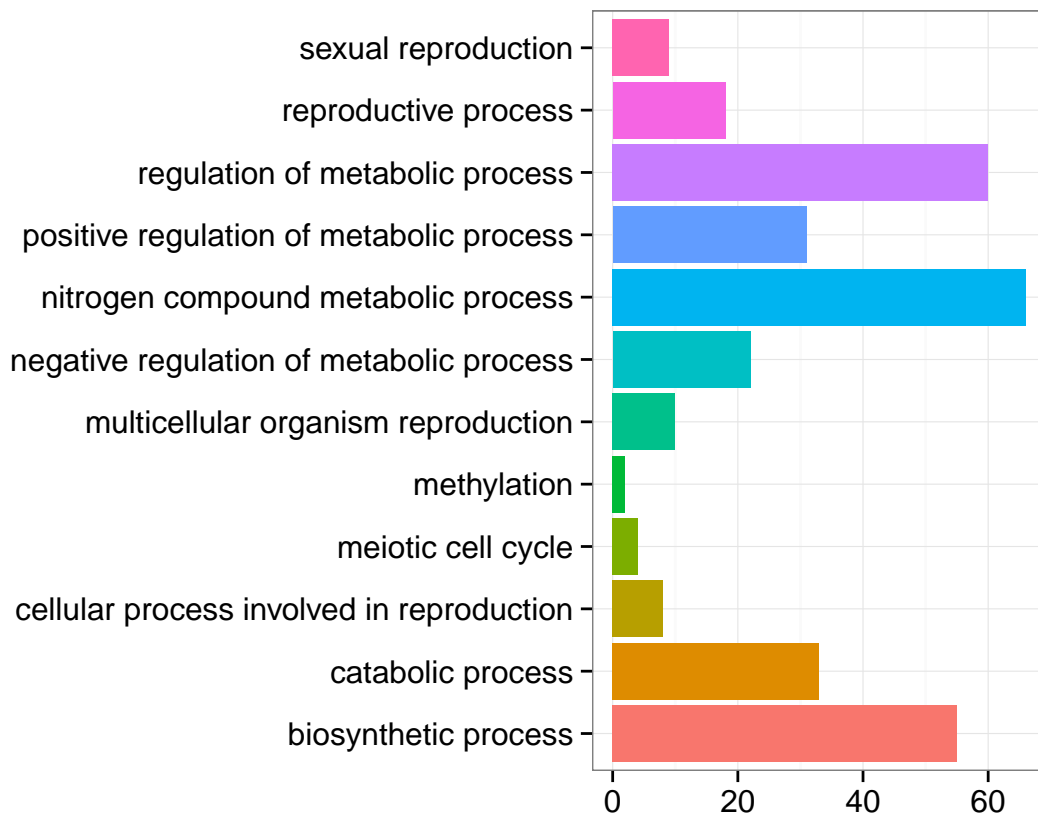
4.7 Visualization

The output of `groupGO`, `enrichGO` and `enrichKEGG` can be visualized by bar plot and category-gene-network plot. It is very common to visualize the enrichment result

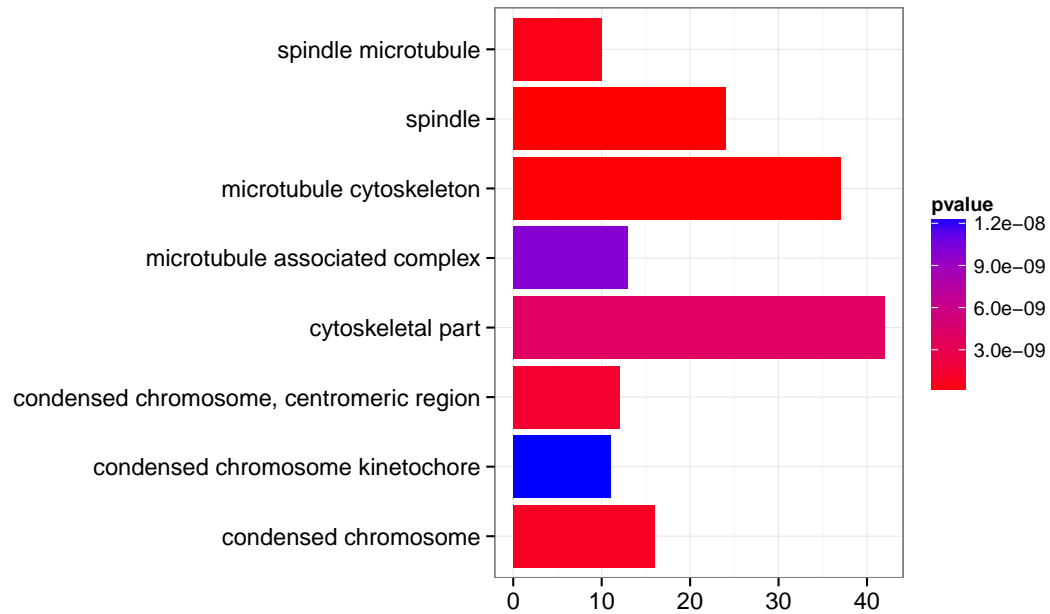
in bar or pie chart. We believe the pie chart is misleading and only provide bar chart.

4.7.1 barplot

```
barplot(ggo, drop = TRUE, showCategory = 12)
```



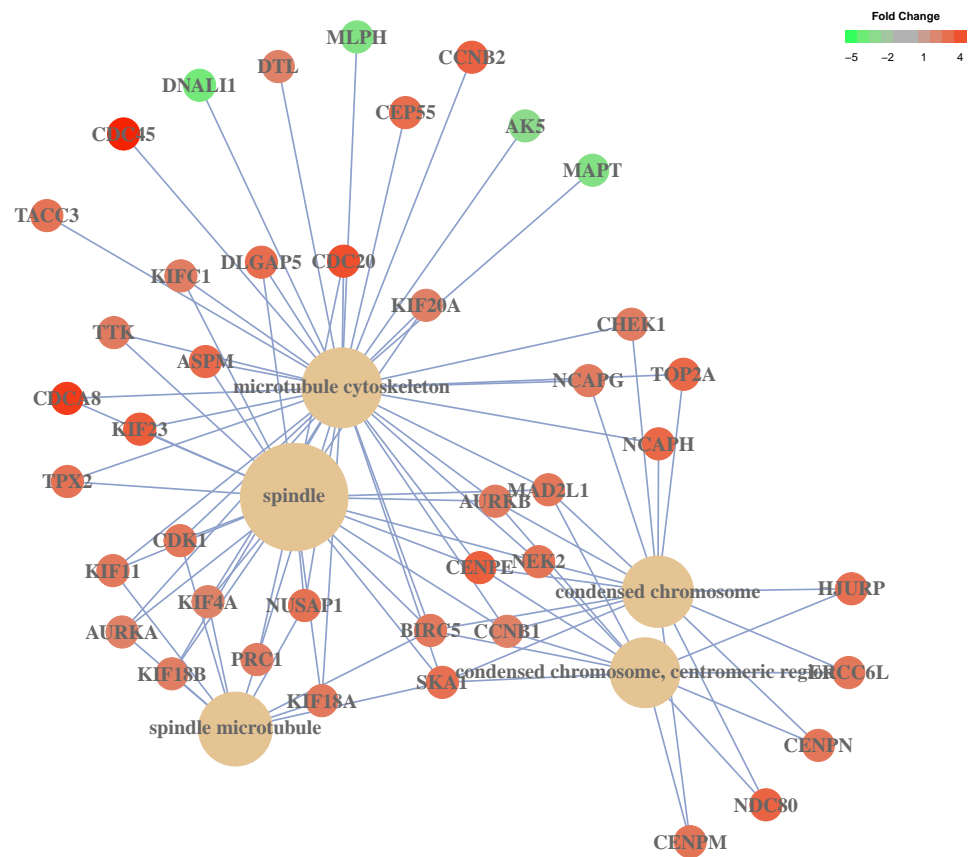
```
barplot(ego, showCategory = 8)
```



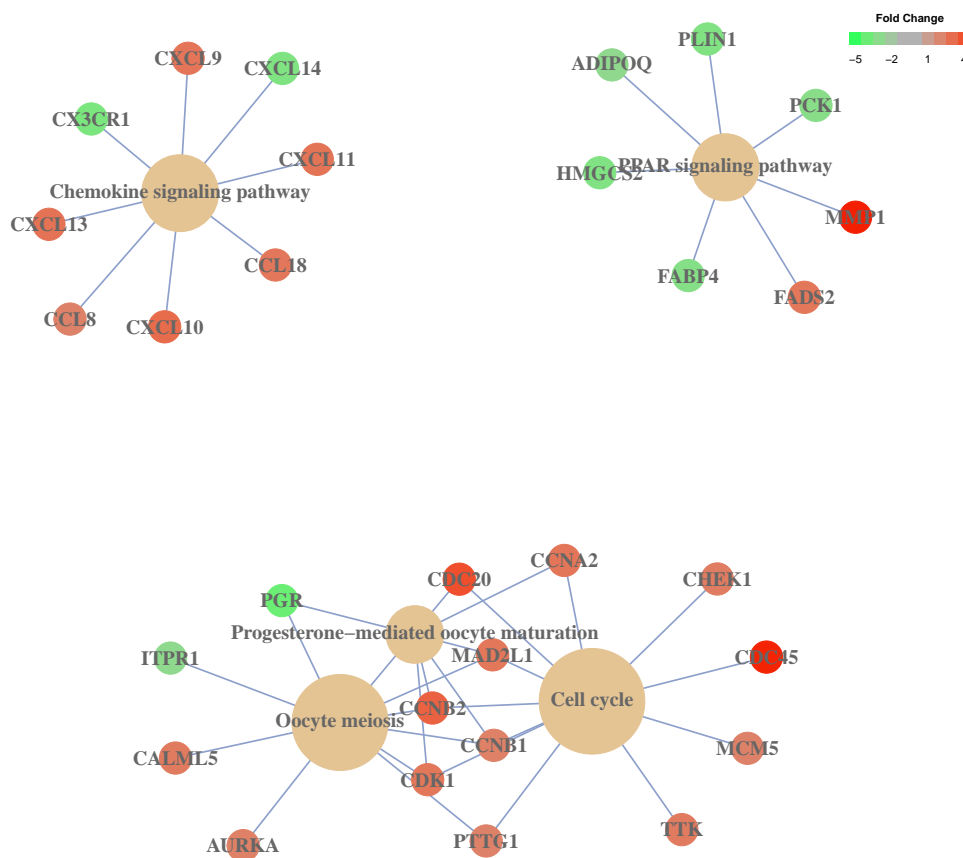
4.7.2 cnetplot

In order to consider the potentially biological complexities in which a gene may belong to multiple annotation categories and provide information of numeric changes if available, we developed `cnetplot` function to extract the complex association.

```
cnetplot(ego, categorySize = "pvalue", foldChange = geneList)
```



```
cnetplot(kk, categorySize = "geneNum", foldChange = geneList)
```

4.7.3 pathview from pathview package

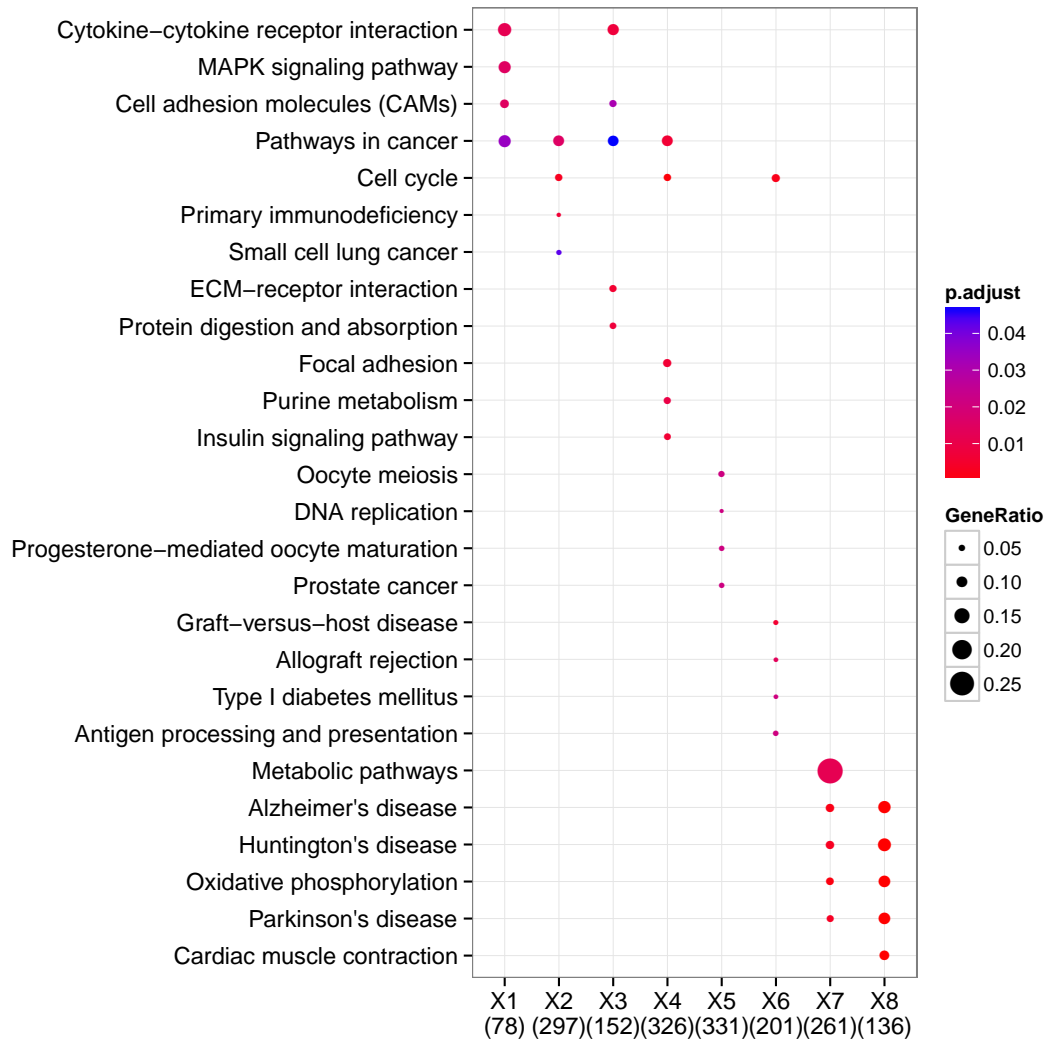
clusterProfiler users can also use *pathview* from the *pathview* [4] to visualize KEGG pathway.

The following example illustrate how to visualize "hsa04110" pathway, which was enriched in our previous analysis.

```
require(pathview)
hsa04110 <- pathview(gene.data = geneList, pathway.id = "hsa04110",
  species = "hsa", limit = list(gene = max(abs(geneList)),
    cpd = 1))

## [1] "Downloading xml files for hsa04110, 1/1 pathways.."
## [1] "Downloading png files for hsa04110, 1/1 pathways.."

## Working in directory /private/tmp/RtmpzyvpDt/Rbuild17c9f1a161697/clusterProfil
## Writing image file hsa04110.pathview.png
```

By default, only top 5 (most significant) categories of each cluster was plotted. User can changes the parameter *showCategory* to specify how many categories of each cluster to be plotted, and if *showCategory* was set to *NULL*, the whole result will be plotted.

The dot sizes were based on their corresponding row percentage by default, and user can set the parameter *by* to "count" to make the comparison based on gene counts. The parameter *by* can also set to "rowPercentage" to normalize the dot sizes, since some categories may contain a large number of genes, and make the dot sizes of those small categories too small to compare. The default parameter *by* is setting to "geneRatio", which corresponding to the "GeneRatio" column of the output. To provide the full information, we also provide number of identified genes in each category (numbers in parentheses) when *by* is setting to "rowPercentage" and number of gene clusters in each cluster label (numbers in parentheses) when *by* is setting to "geneRatio", as shown in Figure 3. If the dot sizes were based on "count", the row numbers will not shown.

The p-values indicate that which categories are more likely to have biological meanings. The dots in the plot are color-coded based on their corresponding

p-values. Color gradient ranging from red to blue correspond to in order of increasing p-values. That is, red indicate low p-values (high enrichment), and blue indicate high p-values (low enrichment). P-values and adjusted p-values were filtered out by the threshold giving by parameter *pvalueCutoff*, and FDR can be estimated by *qvalue*.

User can refer to the example in [2]; we analyzed the publicly available expression dataset of breast tumour tissues from 200 patients (GSE11121, Gene Expression Omnibus) [5]. We identified 8 gene clusters from differentially expressed genes, and using `compareCluster` to compare these gene clusters by their enriched biological process.

Another example was shown in [6], we calculated functional similarities among viral miRNAs using method described in [7], and compared significant KEGG pathways regulated by different viruses using `compareCluster`.

The comparison function was designed as a general-package for comparing gene clusters of any kind of ontology associations, not only `groupGO`, `enrichGO`, and `enrichKEGG` this package provided, but also other biological and biomedical ontologies, for instance, `enrichDO` from *DOSE* and `enrichPathway` from *ReactomePA* work fine with `compareCluster` for comparing biological themes in disease and reactome pathway perspective. More details can be found in the vignettes of *DOSE* and *ReactomePA*.

6 Session Information

The version number of R and packages loaded for generating the vignette were:

- R version 3.1.0 RC (2014-04-02 r65358), x86_64-apple-darwin10.8.0
- Locale:
en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
- Base packages: base, datasets, graphics, grDevices, methods, parallel, stats, utils
- Other packages: AnnotationDbi 1.26.0, Biobase 2.24.0, BiocGenerics 0.10.0, clusterProfiler 1.12.0, DBI 0.2-7, DOSE 2.2.0, GenomeInfoDb 1.0.0, ggplot2 0.9.3.1, GO.db 2.14.0, graph 1.42.0, KEGGgraph 1.22.0, knitr 1.5, org.Hs.eg.db 2.14.0, pathview 1.4.0, RSQLite 0.11.4, XML 3.98-1.1
- Loaded via a namespace (and not attached): Biostrings 2.32.0, codetools 0.2-8, colorspace 1.2-4, dichromat 2.0-0, digest 0.6.4, DO.db 2.7, evaluate 0.5.3, formatR 0.10, GOSemSim 1.22.0, grid 3.1.0, gtable 0.1.2, highr 0.3, httr 0.3, igraph 0.7.0, IRanges 1.21.45, KEGG.db 2.14.0, KEGGREST 1.4.0, labeling 0.2, MASS 7.3-31, munsell 0.4.2, plyr 1.8.1, png 0.1-7, proto 0.3-10, qvalue 1.38.0, RColorBrewer 1.0-5, Rcpp 0.11.1, RCurl 1.95-4.1, reshape2 1.2.2,

Rgraphviz 2.8.0, scales 0.2.3, stats4 3.1.0, stringr 0.6.2, tcltk 3.1.0, tools 3.1.0, XVector 0.4.0, zlibbioc 1.10.0

References

- [1] Guangchuang Yu, Fei Li, Yide Qin, Xiaochen Bo, Yibo Wu, and Shengqi Wang. Gosemsim: an r package for measuring semantic similarity among go terms and gene products. *Bioinformatics*, 26(7):976–978, 2010. PMID: 20179076.
- [2] Guangchuang Yu, Le-Gen Wang, Yanyan Han, and Qing-Yu He. clusterpro-filer: an r package for comparing biological themes among gene clusters. *OMICS: A Journal of Integrative Biology*, 16(5):284–287, May 2012.
- [3] Elizabeth I Boyle, Shuai Weng, Jeremy Gollub, Heng Jin, David Botstein, J Michael Cherry, and Gavin Sherlock. GO::TermFinder—open source software for accessing gene ontology information and finding significantly enriched gene ontology terms associated with a list of genes. *Bioinformatics (Oxford, England)*, 20(18):3710–3715, December 2004. PMID: 15297299.
- [4] Weijun Luo and Cory Brouwer. Pathview: an R/Bioconductor package for pathway-based data integration and visualization. 29:1830–1831. PMID: 23740750.
- [5] Marcus Schmidt, Daniel Böhmer, Christian von Thüne, Eric Steiner, Alexander Puhl, Henryk Pilch, Hans-Anton Lehr, Jan G. Hengstler, Heinz Köhl, and Mathias Gehrmann. The humoral immune system has a key prognostic impact in node-negative breast cancer. *Cancer Research*, 68(13):5405–5413, July 2008.
- [6] Guangchuang Yu and Qing-Yu He. Functional similarity analysis of human virus-encoded miRNAs. *Journal of Clinical Bioinformatics*, 1(1):15, May 2011.
- [7] Guangchuang Yu, Chuan-Le Xiao, Xiaochen Bo, Chun-Hua Lu, Yide Qin, Sheng Zhan, and Qing-Yu He. A new method for measuring functional similarity of microRNAs. *Journal of Integrated OMICS*, 1(1):49–54, February 2011.