

Pathway analysis and Drug repositioning for Psoriasis based on cogen

Zhilong Jia

2015-09-04

Contents

1	Introduction	1
2	Downloading the raw data of GSE13355	1
3	Differential Expression Analysis	2
4	Pathway Analysis by cogen	3
4.1	The heatmap with co-expression information	4
4.2	Figure 2: The result of pathway analysis	4
4.3	Table S1: Make Input for GSEA	4
5	Drug repositioning by cogen	7
5.1	Figure 3: Drug repositioning for cluster 3 (A)	7
5.2	Figure 4: Drug repositioning for cluster 4 (B)	7
5.3	Figure S1: Drug repositioning for cluster 6 (C)	7
5.4	Table S2: Output DEGs for CMAP and NFFinder Analysis	11
6	System Info	11

1 Introduction

*As a reproducible research, this is the whole code and files necessary for the manuscript, **Drug repositioning and drug mode of action discovering based on co-expressed gene-set enrichment analysis.***

2 Downloading the raw data of [GSE13355](#)

```
# Download file from GEO and untar them if nothing in ../data/GSE13355_RAW
if (length(dir("../data/GSE13355_RAW", all.files=FALSE)) ==0) {

  download.file("http://www.ncbi.nlm.nih.gov/geo/download/?acc=GSE13355&format=file",
               destfile="../data/GSE13355_RAW.tar")
  untar("../data/GSE13355_RAW.tar", exdir="../data/GSE13355_RAW")
}
```

3 Differential Expression Analysis

```
library(GEOquery)
library (affy)

# GSE13355
GSE13355raw <- ReadAffy(cefile.path="../data/GSE13355_RAW")
sampleNames(GSE13355raw) <- sub("(_|\\\\.)*CEL\\.gz","", sampleNames(GSE13355raw))

# Sample Label preprocessing
GSE13355series <- getGEO("GSE13355", destdir="../data")
GSE13355label <- pData(GSE13355series$GSE13355_series_matrix.txt.gz)[,c("title", "geo_accession")]
GSE13355label$title <- as.character(GSE13355label$title)

GSE13355label <- GSE13355label[grep("NN", GSE13355label$title, invert = T),]
GSE13355label[grep("PN", GSE13355label$title),"state"] = "ct"
GSE13355label[grep("PP", GSE13355label$title),"state"] = "Psoriasis"
GSE13355label$state <- as.factor(GSE13355label$state)
GSE13355label[, "gse_id"] = "GSE13355"
GSE13355label$rep <- sapply(strsplit(GSE13355label$title, "_"), "[", 2)

GSE13355raw <- GSE13355raw[,as.character(GSE13355label$geo_accession)]

vmd = data.frame(labelDescription = c("title", "geo_accession", "state", "gse_id", "rep"))
phenoData(GSE13355raw) = new("AnnotatedDataFrame", data = GSE13355label, varMetadata = vmd)
pData(protocolData(GSE13355raw)) <-
  pData(protocolData(GSE13355raw))[rownames(GSE13355label),,drop=FALSE]

# RMA normalization
GSE13355rma <- rma(GSE13355raw)

## Background correcting
## Normalizing
## Calculating Expression

#####
# Filter the non-informative and non-expressed genes first.
library(MetaDE)
library(annotate)
library(hgu133plus2.db)

GSE13355.Explist <- list(GSE13355=list(x = exprs(GSE13355rma),
      y = ifelse (GSE13355label$state=="ct", 0, 1),
      symbol = getSYMBOL(rownames(exprs(GSE13355rma)), "hgu133plus2") ))
GSE13355.Explist <- MetaDE.match(GSE13355.Explist, pool.replicate="IQR")
GSE13355.Explist.filtered <- MetaDE.filter(GSE13355.Explist, c(0.2,0.2))
colnames(GSE13355.Explist.filtered$GSE13355$x) <- colnames(exprs(GSE13355rma))

# DEG analysis via limma
DElimma <- function (Expdata, Exlabel){

  library(limma)
```

```

Expdesign <- model.matrix(~as.factor(Explabel$rep) + Explabel$state)
Expfit1 <- lmFit(Expdata, Expdesign)
Expfit2 <- eBayes(Expfit1)
dif_Exp <- topTable(Expfit2, coef=tail(colnames(Expdesign), 1), number=Inf)

return (dif_Exp)
}

GSE13355.limma <- DElimma(GSE13355.Explist.filtered$GSE13355$x, GSE13355label)
GSE13355.DE <- GSE13355.limma[GSE13355.limma$adj.P.Val<=0.05 & abs(GSE13355.limma$logFC)>=1,]
GSE13355.DEG <- rownames(GSE13355.DE)
GSE13355.DEG.expr <- GSE13355.Explist.filtered$GSE13355$x[GSE13355.DEG,]

```

4 Pathway Analysis by cogena

```

# Install cogena if not
# devtools::install_github("zhilongjia/cogena")
library(cogena)
annoGMT <- "c2.cp.kegg.v5.0.symbols.gmt.xz"
annofile <- system.file("extdata", annoGMT, package="cogena")
# nClust <- 2:20
# clMethods <- c("hierarchical", "kmeans", "diana", "fanny", "som", "sota", "pam", "clara", "agnes")
nClust <- 7
clMethods <- c("pam")
ncore <- 7
# Co-expression analysis
genecl_result <- coExp(GSE13355.DEG.expr, nClust=nClust,
                      clMethods=clMethods,
                      metric="correlation",
                      method="complete",
                      ncore=ncore,
                      verbose=FALSE)
sampleLabel <- GSE13355label$state
names(sampleLabel) <- rownames(GSE13355label)
# cogena analysis (Pathway analysis)
cogena_result <- clEnrich(genecl_result, annofile=annofile, sampleLabel=sampleLabel)

# Summary the results obtained by cogena
summary(cogena_result)

```

```

##
## Clustering Methods:
## pam
##
## The Number of Clusters:
## 7
##
## Metric of Distance Matrix:
## correlation
##

```

```
## Agglomeration method for hierarchical clustering (hclust and agnes):
## complete
##
## Gene set:
## c2.cp.kegg.v5.0.symbols.gmt.xz
```

4.1 The heatmap with co-expression information

```
heatmapCluster(cogena_result, "pam", "7", maintitle="Psoriasis")
```

```
## The number of genes in each cluster:
## upDownGene
## 1 2
## 468 238
## cluster_size
## 1 2 3 4 5 6 7
## 257 65 81 130 94 61 18
```

4.2 Figure 2: The result of pathway analysis

```
heatmapPEI(cogena_result, "pam", "7", printGS=FALSE, maintitle="Psoriasis")
```

4.3 Table S1: Make Input for GSEA

This is to get the *Table S1*. The result can be obtained from `../result/GSEA_output`, too. See `gct` and `cls` file format if needed.

```
expData <- as.data.frame(exprs(GSE13355rma))
expData$DESCRIPTION <- NA
expData <- expData[,c("DESCRIPTION", colnames(expData)[1:116])]

write.table(expData, file="../result/GSEA_input/GSE13355_exp.gct", sep="\t", quote=FALSE)
#####
# Add the following 3 lines at the beginning of GSE13355_exp.gct manually.
#1.2
54675 116
NAME
#####
write.table(t(as.character(GSE13355label$state)),
  file="../result/GSEA_input/GSE13355.cls", quote=FALSE, col.names=FALSE,
  row.names=FALSE)
#####
# Add the following 2 lines at the beginning of GSE13355.cls manually
116 2 1
#ct Psoriasis
#####
# Download gsea2-2.1.0.jar from the GSEA website
```

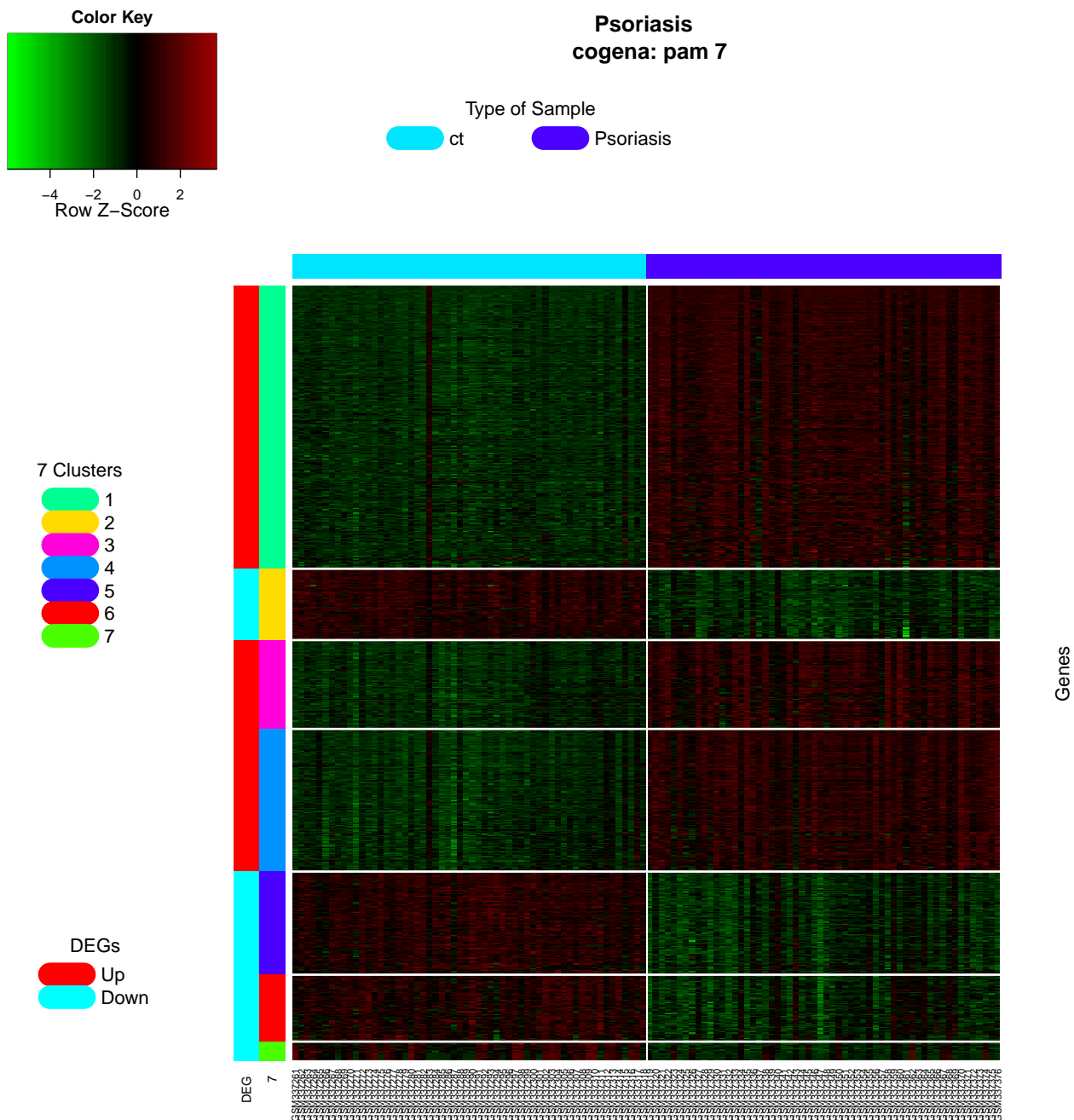


Figure 1: Heatmap with co-expression information

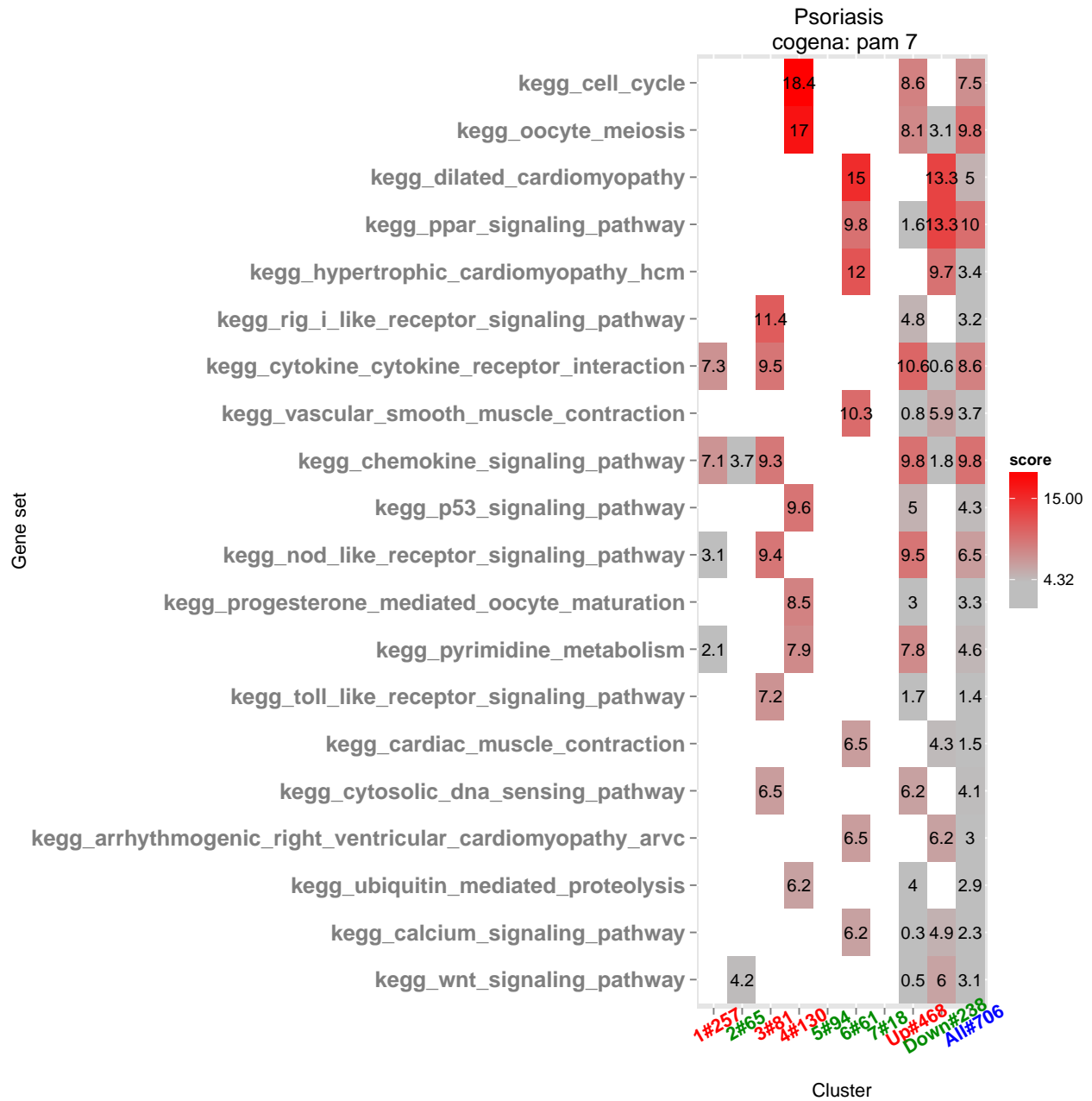


Figure 2: Pathway Analysis

```
# Or from [my github](https://github.com/zhilongjia/geneRanking/blob/master/src/gsea2-2.1.0.jar)
# to the current directory.
#
# GSEA analysis
# java -cp ./gsea2-2.1.0.jar -Xmx512m xtools.gsea.Gsea -res ../result/GSEA_input/GSE13355_exp.gct
# -cls ../result/GSEA_input/GSE13355.cls -gmw ../result/GSEA_input/c2.cp.kegg.v5.0.symbols.gmt
# -collapse true -mode Max_probe -norm meandiv -nperm 1000 -permute phenotype
# -rnd_type no_balance -scoring_scheme weighted -rpt_label GSE13355 -metric Signal2Noise
# -sort real -order descending -chip ../result/GSEA_input/HG_U133_Plus_2.chip
# -include_only_symbols true -make_sets true -median false -num 100 -plot_top_x 20
# -rnd_seed timestamp -save_rnd_lists false -set_max 500 -set_min 15 -zip_report false
# -out ../result/GSEA_output -gui false
```

5 Drug repositioning by cogena

```
# Drug repositioning based on CmapDn100 gene set
cmapDn100_cogena_result <- clEnrich_one(genecl_result, "pam", "7",
  annofile=system.file("extdata", "CmapDn100.gmt.xz", package="cogena"),
  sampleLabel=sampleLabel)

# Drug repositioning based on CmapUp100 gene set
cmapUp100_cogena_result <- clEnrich_one(genecl_result, method="pam", nCluster="7",
  annofile=system.file("extdata", "CmapUp100.gmt.xz", package="cogena"),
  sampleLabel=sampleLabel)
```

5.1 Figure 3: Drug repositioning for cluster 3 (A)

```
heatmapPEI(cmapDn100_cogena_result, "pam", "7", printGS=FALSE,
  orderMethod = "3", maintitle="Psoriasis")
```

5.2 Figure 4: Drug repositioning for cluster 4 (B)

```
heatmapPEI(cmapDn100_cogena_result, "pam", "7", printGS=FALSE,
  orderMethod = "4", maintitle="Psoriasis")
```

5.3 Figure S1: Drug repositioning for cluster 6 (C)

```
# See Figure 5
heatmapPEI(cmapUp100_cogena_result, "pam", "7", printGS=FALSE,
  orderMethod = "6", maintitle="Psoriasis")
```

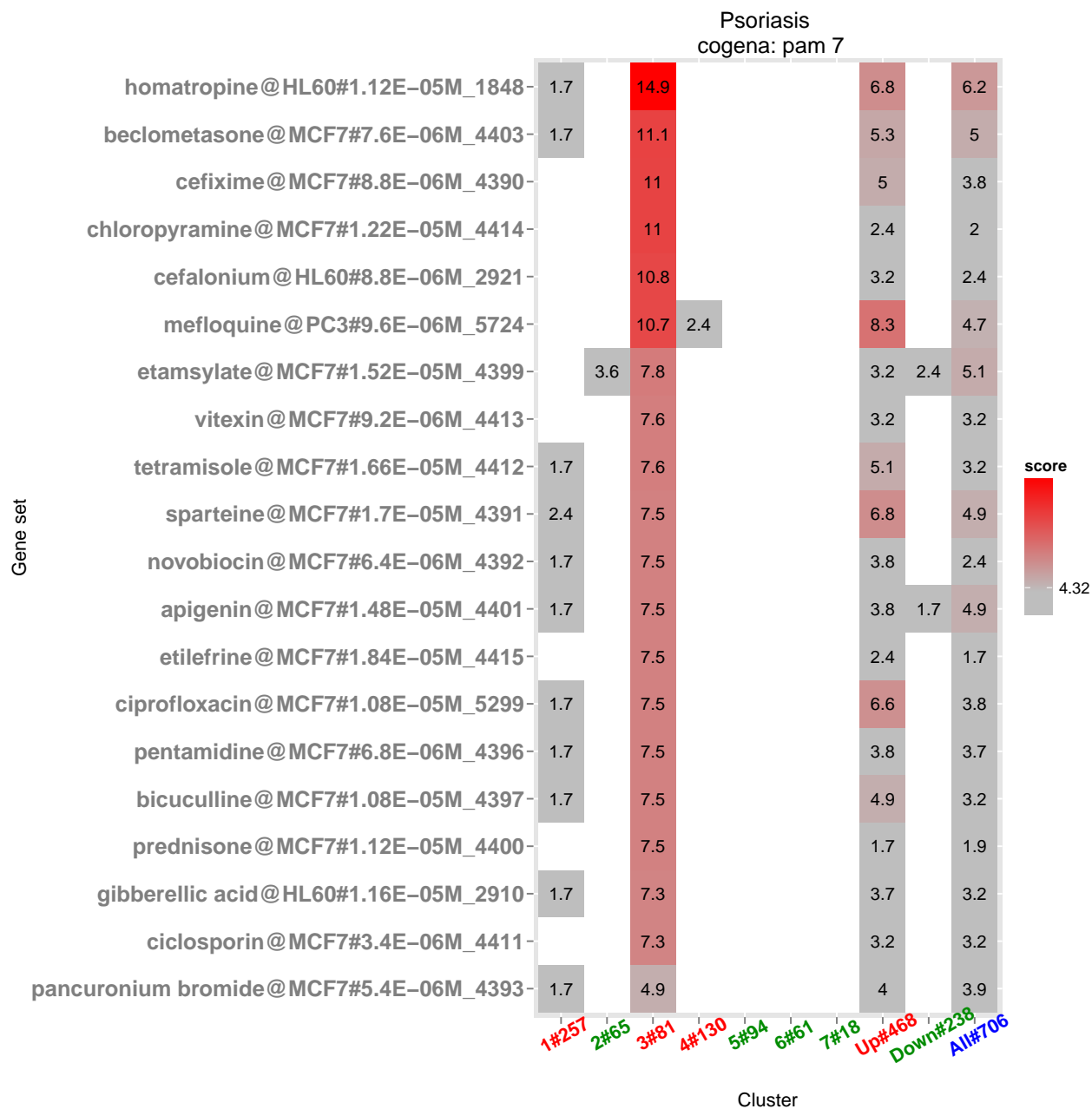


Figure 3: Drug Repositioning for cluster 3

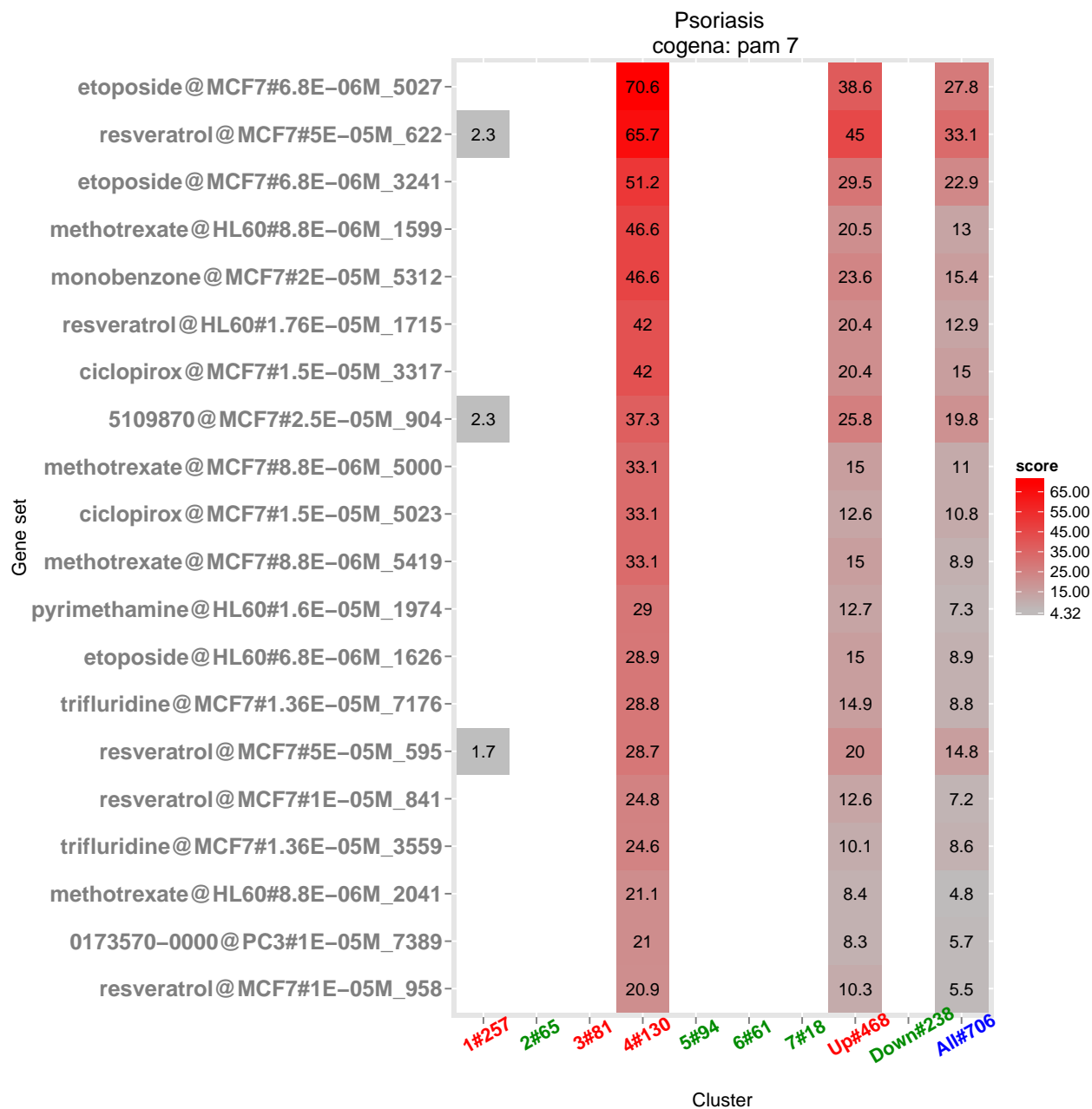


Figure 4: Drug Repositioning for cluster 4

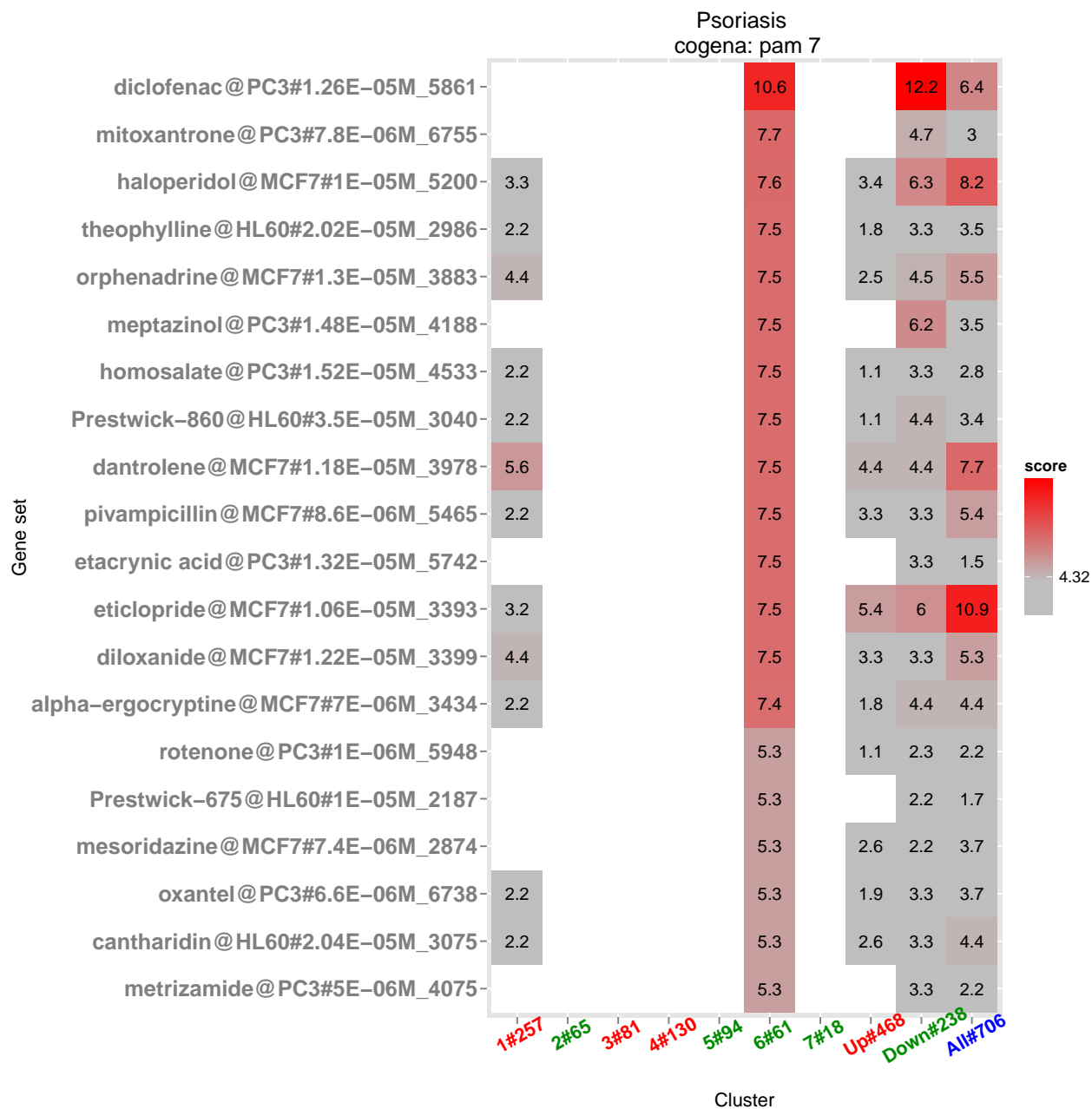


Figure 5: Drug Repositioning for cluster 6

5.4 Table S2: Output DEGs for CMAP and NFFinder Analysis

These outputs are used for [CMap](#) and [NFFinder](#) analysis to get the *Table S2*. The results can be obtained from `../result/CMAP_output` and `../result/NFFinder_output`, as well. For NFFinder, the CMap database and Profile matching, “Inverse”, are used.

```
# Convert gene symbols to probes in HGU133a.
symbol2Probe <- function(gs){
  library(hgu133a.db)
  p <- AnnotationDbi::select(hgu133a.db, gs, "PROBEID", "SYMBOL")$PROBEID
  p <- unique(p[which(!is.na(p))])
}

upGene <- rownames(GSE13355.limma[GSE13355.limma$logFC>= 1 & GSE13355.limma$adj.P.Val<=0.05,])
dnGene <- rownames(GSE13355.limma[GSE13355.limma$logFC<= -1 & GSE13355.limma$adj.P.Val<=0.05,])
upProbe <- symbol2Probe(upGene)
dnProbe <- symbol2Probe(dnGene)

# Output files for CMap and NFFinder
write.table(upProbe, file=paste0("../result/CMAP_input/", "GSE13355_Up.grp"),
  quote=F, col.names = F, row.names = F)
write.table(dnProbe, file=paste0("../result/CMAP_input/", "GSE13355_Dn.grp"),
  quote=F, col.names = F, row.names = F)
write.table(upGene, file=paste0("../result/NFFinder_input/", "GSE13355_Up.txt"),
  quote=F, col.names = F, row.names = F)
write.table(dnGene, file=paste0("../result/NFFinder_input/", "GSE13355_Dn.txt"),
  quote=F, col.names = F, row.names = F)
#####
```

6 System Info

```
## R version 3.2.0 (2015-04-16)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Debian GNU/Linux jessie/sid
##
## locale:
##  [1] LC_CTYPE=en_GB.UTF-8      LC_NUMERIC=C
##  [3] LC_TIME=en_GB.UTF-8      LC_COLLATE=en_GB.UTF-8
##  [5] LC_MONETARY=en_GB.UTF-8  LC_MESSAGES=en_GB.UTF-8
##  [7] LC_PAPER=en_GB.UTF-8     LC_NAME=C
##  [9] LC_ADDRESS=C             LC_TELEPHONE=C
## [11] LC_MEASUREMENT=en_GB.UTF-8 LC_IDENTIFICATION=C
##
## attached base packages:
##  [1] stats4      tools      parallel    stats      graphics  grDevices  utils
##  [8] datasets    methods    base
##
## other attached packages:
##  [1] limma_3.24.14      hgu133plus2.db_3.1.3  org.Hs.eg.db_3.1.2
##  [4] RSQLite_1.0.0      DBI_0.3.1             annotate_1.46.1
##  [7] XML_3.98-1.3       AnnotationDbi_1.30.1  GenomeInfoDb_1.4.1
## [10] IRanges_2.2.5      S4Vectors_0.6.3      MetaDE_1.0.5
```

```

## [13] combinat_0.0-8      impute_1.42.0      survival_2.38-3
## [16] hgu133plus2cdf_2.16.0 affy_1.46.1      GEOquery_2.35.4
## [19] Biobase_2.28.0      BiocGenerics_0.14.0 cogen_1.2.0
## [22] kohonen_2.0.18      MASS_7.3-43       class_7.3-13
## [25] ggplot2_1.0.1      cluster_2.0.3
##
## loaded via a namespace (and not attached):
## [1] splines_3.2.0      foreach_1.4.2      gtools_3.5.0
## [4] assertthat_0.1     amap_0.8-14        yaml_2.1.13
## [7] robustbase_0.92-5  corrplot_0.73      lattice_0.20-33
## [10] digest_0.6.8       colorspace_1.2-6   htmltools_0.2.6
## [13] preprocessCore_1.30.0 Matrix_1.2-2       plyr_1.8.3
## [16] pcaPP_1.9-60       devtools_1.8.0     zlibbioc_1.14.0
## [19] xtable_1.7-4       mvtnorm_1.0-3      scales_0.2.5
## [22] gdata_2.17.0       affyio_1.36.0      git2r_0.10.1
## [25] biwt_1.0           fastcluster_1.1.16 lazyeval_0.1.10
## [28] proto_0.3-10       magrittr_1.5       mclust_5.0.2
## [31] memoise_0.2.1      evaluate_0.7       apcluster_1.4.1
## [34] doParallel_1.0.8   gplots_2.17.0     xml2_0.1.1
## [37] BiocInstaller_1.18.4 formatR_1.2        stringr_1.0.0
## [40] munsell_0.4.2      rversions_1.0.2    compiler_3.2.0
## [43] caTools_1.17.1     grid_3.2.0         RCurl_1.95-4.7
## [46] iterators_1.0.7    bitops_1.0-6       rmarkdown_0.7
## [49] gtable_0.1.2       codetools_0.2-14   curl_0.9.1
## [52] reshape2_1.4.1     rrcov_1.3-8        R6_2.1.0
## [55] knitr_1.10.5       dplyr_0.4.2        KernSmooth_2.23-15
## [58] stringi_0.5-5      Rcpp_0.12.0        DEoptimR_1.0-3

```