

Class14: RNASeq Mini-project

Chelsea A16871799

The authors report on differential analysis of lung fibroblasts in response to loss of the developmental transcription factor HOXA1. Their results and others indicate that HOXA1 is required for lung fibroblast and HeLa cell cycle progression. In particular their analysis show that “loss of HOXA1 results in significant expression level changes in thousands of individual transcripts, along with isoform switching events in key regulators of the cell cycle”. For our session we have used their Sailfish gene-level estimated counts and hence are restricted to protein-coding genes only.

##Data Import

```
metaFile <- "GSE37704_metadata.csv"
countFile <- "GSE37704_featurecounts.csv"

# Import metadata and take a peak
colData = read.csv(metaFile, row.names=1)
head(colData)
```

```
              condition
SRR493366 control_sirna
SRR493367 control_sirna
SRR493368 control_sirna
SRR493369      hoxa1_kd
SRR493370      hoxa1_kd
SRR493371      hoxa1_kd
```

```
countData=read.csv(countFile, row.names=1)
head(countData)
```

```
              length SRR493366 SRR493367 SRR493368 SRR493369 SRR493370
ENSG00000186092      918          0          0          0          0
```

ENSG00000279928	718	0	0	0	0	0
ENSG00000279457	1982	23	28	29	29	28
ENSG00000278566	939	0	0	0	0	0
ENSG00000273547	939	0	0	0	0	0
ENSG00000187634	3214	124	123	205	207	212
	SRR493371					
ENSG00000186092	0					
ENSG00000279928	0					
ENSG00000279457	46					
ENSG00000278566	0					
ENSG00000273547	0					
ENSG00000187634	258					

we need the countData and colData files to match up so we will need to remove that odd first column in countData namely countData\$length

```
countData <- as.matrix(countData[,-1])
head(countData)
```

	SRR493366	SRR493367	SRR493368	SRR493369	SRR493370	SRR493371
ENSG00000186092	0	0	0	0	0	0
ENSG00000279928	0	0	0	0	0	0
ENSG00000279457	23	28	29	29	28	46
ENSG00000278566	0	0	0	0	0	0
ENSG00000273547	0	0	0	0	0	0
ENSG00000187634	124	123	205	207	212	258

Q. Complete the code below to filter countData to exclude genes (i.e. rows) where we have 0 read count across all samples (i.e. columns).

```
to.rm.inds <- rowSums(countData)==0

countData =countData[!to.rm.inds,]
head(countData)
```

	SRR493366	SRR493367	SRR493368	SRR493369	SRR493370	SRR493371
ENSG00000279457	23	28	29	29	28	46
ENSG00000187634	124	123	205	207	212	258
ENSG00000188976	1637	1831	2383	1226	1326	1504
ENSG00000187961	120	153	180	236	255	357

ENSG00000187583	24	48	65	44	48	64
ENSG00000187642	4	9	16	14	16	16

```
nrow(countData)
```

```
[1] 15975
```

```
##DESeq Setup and analysis
```

```
library(DESeq2)
```

```
Loading required package: S4Vectors
```

```
Loading required package: stats4
```

```
Loading required package: BiocGenerics
```

```
Attaching package: 'BiocGenerics'
```

```
The following objects are masked from 'package:stats':
```

```
IQR, mad, sd, var, xtabs
```

```
The following objects are masked from 'package:base':
```

```
anyDuplicated, aperm, append, as.data.frame, basename, cbind,
colnames, dirname, do.call, duplicated, eval, evalq, Filter, Find,
get, grep, grepl, intersect, is.unsorted, lapply, Map, mapply,
match, mget, order, paste, pmax, pmax.int, pmin, pmin.int,
Position, rank, rbind, Reduce, rownames, sapply, setdiff, sort,
table, tapply, union, unique, unsplit, which.max, which.min
```

```
Attaching package: 'S4Vectors'
```

```
The following object is masked from 'package:utils':
```

```
findMatches
```

The following objects are masked from 'package:base':

expand.grid, I, unname

Loading required package: IRanges

Loading required package: GenomicRanges

Loading required package: GenomeInfoDb

Loading required package: SummarizedExperiment

Warning: package 'SummarizedExperiment' was built under R version 4.3.2

Loading required package: MatrixGenerics

Loading required package: matrixStats

Attaching package: 'MatrixGenerics'

The following objects are masked from 'package:matrixStats':

colAlls, colAnyNAs, colAnys, colAvgPerRowSet, colCollapse,
colCounts, colCummaxs, colCummins, colCumprods, colCumsums,
colDiffs, colIQRDiffs, colIQRs, colLogSumExps, colMadDiffs,
colMads, colMaxs, colMeans2, colMedians, colMins, colOrderStats,
colProds, colQuantiles, colRanges, colRanks, colSdDiffs, colSds,
colSums2, colTabulates, colVarDiffs, colVars, colWeightedMads,
colWeightedMeans, colWeightedMedians, colWeightedSds,
colWeightedVars, rowAlls, rowAnyNAs, rowAnys, rowAvgPerColSet,
rowCollapse, rowCounts, rowCummaxs, rowCummins, rowCumprods,
rowCumsums, rowDiffs, rowIQRDiffs, rowIQRs, rowLogSumExps,
rowMadDiffs, rowMads, rowMaxs, rowMeans2, rowMedians, rowMins,
rowOrderStats, rowProds, rowQuantiles, rowRanges, rowRanks,
rowSdDiffs, rowSds, rowSums2, rowTabulates, rowVarDiffs, rowVars,
rowWeightedMads, rowWeightedMeans, rowWeightedMedians,
rowWeightedSds, rowWeightedVars

Loading required package: Biobase

Welcome to Bioconductor

Vignettes contain introductory material; view with
'browseVignettes()'. To cite Bioconductor, see
'citation("Biobase")', and for packages 'citation("pkgname")'.

Attaching package: 'Biobase'

The following object is masked from 'package:MatrixGenerics':

rowMedians

The following objects are masked from 'package:matrixStats':

anyMissing, rowMedians

```
dds = DESeqDataSetFromMatrix(countData=countData,  
                              colData=colData,  
                              design=~condition)
```

Warning in DESeqDataSet(se, design = design, ignoreRank): some variables in
design formula are characters, converting to factors

```
dds = DESeq(dds)
```

estimating size factors

estimating dispersions

gene-wise dispersion estimates

mean-dispersion relationship

final dispersion estimates

fitting model and testing

```
dds
```

```
class: DESeqDataSet
dim: 15975 6
metadata(1): version
assays(4): counts mu H cooks
rownames(15975): ENSG00000279457 ENSG00000187634 ... ENSG00000276345
               ENSG00000271254
rowData names(22): baseMean baseVar ... deviance maxCooks
colnames(6): SRR493366 SRR493367 ... SRR493370 SRR493371
colData names(2): condition sizeFactor
```

```
res = results(dds, contrast=c("condition", "hoxa1_kd", "control_sirna"))
```

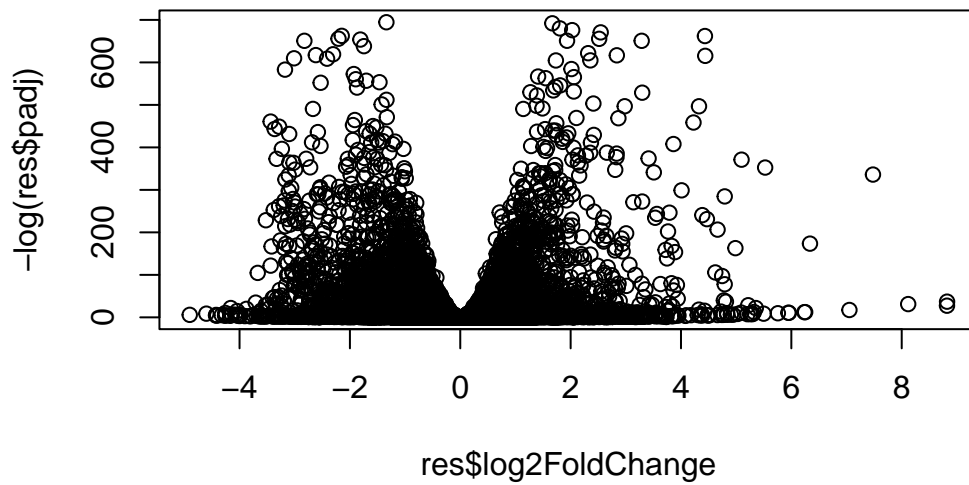
Q. Call the `summary()` function on your results to get a sense of how many genes are up or down-regulated at the default 0.1 p-value cutoff.

```
summary(res)
```

```
out of 15975 with nonzero total read count
adjusted p-value < 0.1
LFC > 0 (up)      : 4349, 27%
LFC < 0 (down)    : 4396, 28%
outliers [1]      : 0, 0%
low counts [2]    : 1237, 7.7%
(mean count < 0)
[1] see 'cooksCutoff' argument of ?results
[2] see 'independentFiltering' argument of ?results
```

volcano plot

```
plot( res$log2FoldChange, -log(res$padj) )
```



Q. Improve this plot by completing the below code, which adds color and axis labels

```
# Make a color vector for all genes
mycols <- rep("gray", nrow(res) )

# Color red the genes with absolute fold change above 2
mycols[ abs(res$log2FoldChange) > 2 ] <- "red"

# Color blue those with adjusted p-value less than 0.01
# and absolute fold change more than 2
inds <- (res$padj<0.01) & (abs(res$log2FoldChange) > 2 )
mycols[ inds ] <- "blue"

plot( res$log2FoldChange, -log(res$padj), col=mycols, xlab="Log2(FoldChange)", ylab="-Log(
```



Q. Use the `mapIDs()` function multiple times to add SYMBOL, ENTREZID and GENENAME annotation to our results by completing the code below.

```
library("AnnotationDbi")
library("org.Hs.eg.db")
```

```
columns(org.Hs.eg.db)
```

```
[1] "ACCNUM"      "ALIAS"       "ENSEMBL"     "ENSEMBLPROT" "ENSEMBLTRANS"
[6] "ENTREZID"    "ENZYME"      "EVIDENCE"     "EVIDENCEALL"  "GENENAME"
[11] "GENETYPE"    "GO"          "GOALL"        "IPI"          "MAP"
[16] "OMIM"        "ONTOLOGY"    "ONTOLOGYALL" "PATH"         "PFAM"
[21] "PMID"        "PROSITE"     "REFSEQ"       "SYMBOL"       "UCSCKG"
[26] "UNIPROT"
```

```
res$symbol = mapIds(org.Hs.eg.db,
                    keys=row.names(res),
                    keytype="ENSEMBL",
```



```
column="SYMBOL",
multiVals="first")
```

'select()' returned 1:many mapping between keys and columns

```
res$entrez = mapIds(org.Hs.eg.db,
  keys=row.names(res),
  keytype="ENSEMBL",
  column="ENTREZID",
  multiVals="first")
```

'select()' returned 1:many mapping between keys and columns

```
res$name = mapIds(org.Hs.eg.db,
  keys=row.names(res),
  keytype="ENSEMBL",
  column="GENENAME",
  multiVals="first")
```

'select()' returned 1:many mapping between keys and columns

```
head(res, 10)
```

log2 fold change (MLE): condition hoxa1_kd vs control_sirna

Wald test p-value: condition hoxa1 kd vs control sirna

DataFrame with 10 rows and 9 columns

	baseMean	log2FoldChange	lfcSE	stat	pvalue
	<numeric>	<numeric>	<numeric>	<numeric>	<numeric>
ENSG00000279457	29.913579	0.1792571	0.3248216	0.551863	5.81042e-01
ENSG00000187634	183.229650	0.4264571	0.1402658	3.040350	2.36304e-03
ENSG00000188976	1651.188076	-0.6927205	0.0548465	-12.630158	1.43989e-36
ENSG00000187961	209.637938	0.7297556	0.1318599	5.534326	3.12428e-08
ENSG00000187583	47.255123	0.0405765	0.2718928	0.149237	8.81366e-01
ENSG00000187642	11.979750	0.5428105	0.5215599	1.040744	2.97994e-01
ENSG00000188290	108.922128	2.0570638	0.1969053	10.446970	1.51282e-25
ENSG00000187608	350.716868	0.2573837	0.1027266	2.505522	1.22271e-02
ENSG00000188157	9128.439422	0.3899088	0.0467163	8.346304	7.04321e-17

	padj	symbol	entrez	name
	<numeric>	<character>	<character>	<character>
ENSG00000237330	0.158192	0.7859552	4.0804729	0.192614 8.47261e-01
ENSG00000279457	6.86555e-01	NA	NA	NA
ENSG00000187634	5.15718e-03	SAMD11	148398	sterile alpha motif ..
ENSG00000188976	1.76549e-35	NOC2L	26155	NOC2 like nucleolar ..
ENSG00000187961	1.13413e-07	KLHL17	339451	kelch like family me..
ENSG00000187583	9.19031e-01	PLEKHN1	84069	pleckstrin homology ..
ENSG00000187642	4.03379e-01	PERM1	84808	PPARGC1 and ESRR ind..
ENSG00000188290	1.30538e-24	HES4	57801	hes family bHLH tran..
ENSG00000187608	2.37452e-02	ISG15	9636	ISG15 ubiquitin like..
ENSG00000188157	4.21963e-16	AGRN	375790	agrin
ENSG00000237330	NA	RNF223	401934	ring finger protein ..

Save results

Q. Finally for this section let's reorder these results by adjusted p-value and save them to a CSV file in your current project directory.

```
res = res[order(res$pvalue),]
write.csv(res, file="deseq_results.csv")
```

Geneset enrichment Now we can load the packages and setup the KEGG data-sets we need

```
library(pathview)
```

```
#####
Pathview is an open source software package distributed under GNU General
Public License version 3 (GPLv3). Details of GPLv3 is available at
http://www.gnu.org/licenses/gpl-3.0.html. Particullary, users are required to
formally cite the original Pathview paper (not just mention it) in publications
or products. For details, do citation("pathview") within R.
```

The pathview downloads and uses KEGG data. Non-academic uses may require a KEGG license agreement (details at <http://www.kegg.jp/kegg/legal.html>).

```
#####
```

```
library(gage)
```

```

library(gageData)

data(kegg.sets.hs)
data(sigmet.idx.hs)

# Focus on signaling and metabolic pathways only
kegg.sets.hs = kegg.sets.hs[sigmet.idx.hs]

# Examine the first 3 pathways
head(kegg.sets.hs, 3)

$`hsa00232 Caffeine metabolism`
[1] "10"    "1544" "1548" "1549" "1553" "7498" "9"

$`hsa00983 Drug metabolism - other enzymes`
[1] "10"    "1066" "10720" "10941" "151531" "1548" "1549" "1551"
[9] "1553"  "1576" "1577"  "1806" "1807"  "1890" "221223" "2990"
[17] "3251"  "3614" "3615"  "3704" "51733" "54490" "54575" "54576"
[25] "54577" "54578" "54579" "54600" "54657" "54658" "54659" "54963"
[33] "574537" "64816" "7083"  "7084" "7172"  "7363" "7364"  "7365"
[41] "7366"  "7367" "7371"  "7372" "7378"  "7498" "79799" "83549"
[49] "8824"  "8833" "9"      "978"

$`hsa00230 Purine metabolism`
[1] "100"    "10201" "10606" "10621" "10622" "10623" "107"    "10714"
[9] "108"    "10846" "109"    "111"    "11128" "11164" "112"    "113"
[17] "114"    "115"    "122481" "122622" "124583" "132"    "158"    "159"
[25] "1633"    "171568" "1716"    "196883" "203"    "204"    "205"    "221823"
[33] "2272"    "22978" "23649"    "246721" "25885" "2618"    "26289" "270"
[41] "271"    "27115" "272"    "2766"    "2977"    "2982"    "2983"    "2984"
[49] "2986"    "2987" "29922"    "3000"    "30833" "30834" "318"    "3251"
[57] "353"    "3614"    "3615"    "3704"    "377841" "471"    "4830"    "4831"
[65] "4832"    "4833"    "4860"    "4881"    "4882"    "4907"    "50484"    "50940"
[73] "51082"    "51251" "51292"    "5136"    "5137"    "5138"    "5139"    "5140"
[81] "5141"    "5142"    "5143"    "5144"    "5145"    "5146"    "5147"    "5148"
[89] "5149"    "5150"    "5151"    "5152"    "5153"    "5158"    "5167"    "5169"
[97] "51728"    "5198"    "5236"    "5313"    "5315"    "53343"    "54107"    "5422"
[105] "5424"    "5425"    "5426"    "5427"    "5430"    "5431"    "5432"    "5433"
[113] "5434"    "5435"    "5436"    "5437"    "5438"    "5439"    "5440"    "5441"

```

```
[121] "5471"    "548644" "55276"  "5557"   "5558"   "55703"  "55811"  "55821"
[129] "5631"    "5634"    "56655"  "56953"  "56985"  "57804"  "58497"  "6240"
[137] "6241"    "64425"   "646625" "654364" "661"    "7498"   "8382"   "84172"
[145] "84265"   "84284"   "84618"  "8622"   "8654"   "87178"  "8833"   "9060"
[153] "9061"    "93034"   "953"    "9533"   "954"    "955"    "956"    "957"
[161] "9583"    "9615"
```

```
foldchanges = res$log2FoldChange
names(foldchanges) = res$entrez
head(foldchanges)
```

```
      1266      54855      1465      51232      2034      2317
-2.422719  3.201955 -2.313738 -2.059631 -1.888019 -1.649792
```

```
# Get the results
keggres = gage(foldchanges, gsets=kegg.sets.hs)
```

```
attributes(keggres)
```

```
$names
[1] "greater" "less"    "stats"
```

```
# Look at the first few down (less) pathways
head(keggres$less)
```

	p.geomean	stat.mean	p.val
hsa04110 Cell cycle	8.995727e-06	-4.378644	8.995727e-06
hsa03030 DNA replication	9.424076e-05	-3.951803	9.424076e-05
hsa03013 RNA transport	1.375901e-03	-3.028500	1.375901e-03
hsa03440 Homologous recombination	3.066756e-03	-2.852899	3.066756e-03
hsa04114 Oocyte meiosis	3.784520e-03	-2.698128	3.784520e-03
hsa00010 Glycolysis / Gluconeogenesis	8.961413e-03	-2.405398	8.961413e-03

	q.val	set.size	exp1
hsa04110 Cell cycle	0.001448312	121	8.995727e-06
hsa03030 DNA replication	0.007586381	36	9.424076e-05
hsa03013 RNA transport	0.073840037	144	1.375901e-03
hsa03440 Homologous recombination	0.121861535	28	3.066756e-03
hsa04114 Oocyte meiosis	0.121861535	102	3.784520e-03
hsa00010 Glycolysis / Gluconeogenesis	0.212222694	53	8.961413e-03

```
pathview(gene.data=foldchanges, pathway.id="hsa04110")
```

Info: Working in directory /Users/chelseazhong/Desktop/bimm143/class14

The diagram illustrates the cell cycle pathway, showing the progression from G1 to S, G2, and M phases. Key components and interactions include:

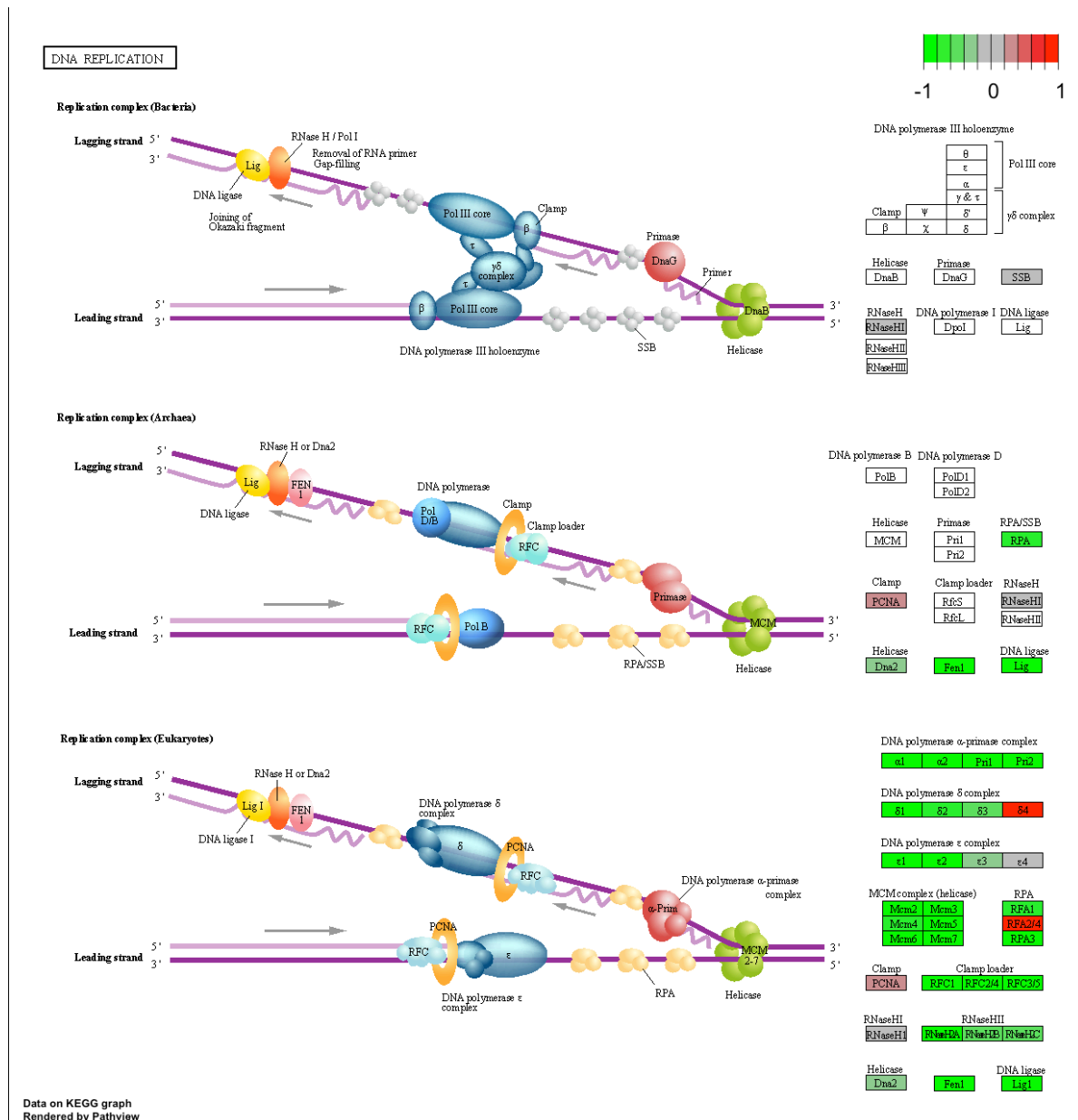
- G1 Phase:** Growth factor withdrawal leads to GSK3β activation. MAPK signaling pathway is involved. Key proteins include p107, E2F4.5, DP-1.2, c-Myc, Miz1, p16, p15, p18, p19, Ink4a, Ink4b, Ink4c, Ink4d, p27.57, p21, SCF, Skp2, TGFβ, Smad2.3, Smad4, ARF, Mdm2, Rb, p300, DNA-PK, ATM/ATR, p53, GADD45, 14-3-3σ, Chk1, 2, PCNA, Cdc25A, Cdc25B, Cdc25C, Cdc2, Cdk1, Cdk2, Cdk4, Cdk6, Cdk7, Cdk8, Cdk9, Cdk10, Cdk11, Cdk12, Cdk13, Cdk14, Cdk15, Cdk16, Cdk17, Cdk18, Cdk19, Cdk20, Cdk21, Cdk22, Cdk23, Cdk24, Cdk25, Cdk26, Cdk27, Cdk28, Cdk29, Cdk30, Cdk31, Cdk32, Cdk33, Cdk34, Cdk35, Cdk36, Cdk37, Cdk38, Cdk39, Cdk40, Cdk41, Cdk42, Cdk43, Cdk44, Cdk45, Cdk46, Cdk47, Cdk48, Cdk49, Cdk50, Cdk51, Cdk52, Cdk53, Cdk54, Cdk55, Cdk56, Cdk57, Cdk58, Cdk59, Cdk60, Cdk61, Cdk62, Cdk63, Cdk64, Cdk65, Cdk66, Cdk67, Cdk68, Cdk69, Cdk70, Cdk71, Cdk72, Cdk73, Cdk74, Cdk75, Cdk76, Cdk77, Cdk78, Cdk79, Cdk80, Cdk81, Cdk82, Cdk83, Cdk84, Cdk85, Cdk86, Cdk87, Cdk88, Cdk89, Cdk90, Cdk91, Cdk92, Cdk93, Cdk94, Cdk95, Cdk96, Cdk97, Cdk98, Cdk99, Cdk100, Cdk101, Cdk102, Cdk103, Cdk104, Cdk105, Cdk106, Cdk107, Cdk108, Cdk109, Cdk110, Cdk111, Cdk112, Cdk113, Cdk114, Cdk115, Cdk116, Cdk117, Cdk118, Cdk119, Cdk120, Cdk121, Cdk122, Cdk123, Cdk124, Cdk125, Cdk126, Cdk127, Cdk128, Cdk129, Cdk130, Cdk131, Cdk132, Cdk133, Cdk134, Cdk135, Cdk136, Cdk137, Cdk138, Cdk139, Cdk140, Cdk141, Cdk142, Cdk143, Cdk144, Cdk145, Cdk146, Cdk147, Cdk148, Cdk149, Cdk150, Cdk151, Cdk152, Cdk153, Cdk154, Cdk155, Cdk156, Cdk157, Cdk158, Cdk159, Cdk160, Cdk161, Cdk162, Cdk163, Cdk164, Cdk165, Cdk166, Cdk167, Cdk168, Cdk169, Cdk170, Cdk171, Cdk172, Cdk173, Cdk174, Cdk175, Cdk176, Cdk177, Cdk178, Cdk179, Cdk180, Cdk181, Cdk182, Cdk183, Cdk184, Cdk185, Cdk186, Cdk187, Cdk188, Cdk189, Cdk190, Cdk191, Cdk192, Cdk193, Cdk194, Cdk195, Cdk196, Cdk197, Cdk198, Cdk199, Cdk200, Cdk201, Cdk202, Cdk203, Cdk204, Cdk205, Cdk206, Cdk207, Cdk208, Cdk209, Cdk210, Cdk211, Cdk212, Cdk213, Cdk214, Cdk215, Cdk216, Cdk217, Cdk218, Cdk219, Cdk220, Cdk221, Cdk222, Cdk223, Cdk224, Cdk225, Cdk226, Cdk227, Cdk228, Cdk229, Cdk230, Cdk231, Cdk232, Cdk233, Cdk234, Cdk235, Cdk236, Cdk237, Cdk238, Cdk239, Cdk240, Cdk241, Cdk242, Cdk243, Cdk244, Cdk245, Cdk246, Cdk247, Cdk248, Cdk249, Cdk250, Cdk251, Cdk252, Cdk253, Cdk254, Cdk255, Cdk256, Cdk257, Cdk258, Cdk259, Cdk260, Cdk261, Cdk262, Cdk263, Cdk264, Cdk265, Cdk266, Cdk267, Cdk268, Cdk269, Cdk270, Cdk271, Cdk272, Cdk273, Cdk274, Cdk275, Cdk276, Cdk277, Cdk278, Cdk279, Cdk280, Cdk281, Cdk282, Cdk283, Cdk284, Cdk285, Cdk286, Cdk287, Cdk288, Cdk289, Cdk290, Cdk291, Cdk292, Cdk293, Cdk294, Cdk295, Cdk296, Cdk297, Cdk298, Cdk299, Cdk300, Cdk301, Cdk302, Cdk303, Cdk304, Cdk305, Cdk306, Cdk307, Cdk308, Cdk309, Cdk310, Cdk311, Cdk312, Cdk313, Cdk314, Cdk315, Cdk316, Cdk317, Cdk318, Cdk319, Cdk320, Cdk321, Cdk322, Cdk323, Cdk324, Cdk325, Cdk326, Cdk327, Cdk328, Cdk329, Cdk330, Cdk331, Cdk332, Cdk333, Cdk334, Cdk335, Cdk336, Cdk337, Cdk338, Cdk339, Cdk340, Cdk341, Cdk342, Cdk343, Cdk344, Cdk345, Cdk346, Cdk347, Cdk348, Cdk349, Cdk350, Cdk351, Cdk352, Cdk353, Cdk354, Cdk355, Cdk356, Cdk357, Cdk358, Cdk359, Cdk360, Cdk361, Cdk362, Cdk363, Cdk364, Cdk365, Cdk366, Cdk367, Cdk368, Cdk369, Cdk370, Cdk371, Cdk372, Cdk373, Cdk374, Cdk375, Cdk376, Cdk377, Cdk378, Cdk379, Cdk380, Cdk381, Cdk382, Cdk383, Cdk384, Cdk385, Cdk386, Cdk387, Cdk388, Cdk389, Cdk390, Cdk391, Cdk392, Cdk393, Cdk394, Cdk395, Cdk396, Cdk397, Cdk398, Cdk399, Cdk400, Cdk401, Cdk402, Cdk403, Cdk404, Cdk405, Cdk406, Cdk407, Cdk408, Cdk409, Cdk410, Cdk411, Cdk412, Cdk413, Cdk414, Cdk415, Cdk416, Cdk417, Cdk418, Cdk419, Cdk420, Cdk421, Cdk422, Cdk423, Cdk424, Cdk425, Cdk426, Cdk427, Cdk428, Cdk429, Cdk430, Cdk431, Cdk432, Cdk433, Cdk434, Cdk435, Cdk436, Cdk437, Cdk438, Cdk439, Cdk440, Cdk441, Cdk442, Cdk443, Cdk444, Cdk445, Cdk446, Cdk447, Cdk448, Cdk449, Cdk450, Cdk451, Cdk452, Cdk453, Cdk454, Cdk455, Cdk456, Cdk457, Cdk458, Cdk459, Cdk460, Cdk461, Cdk462, Cdk463, Cdk464, Cdk465, Cdk466, Cdk467, Cdk468, Cdk469, Cdk470, Cdk471, Cdk472, Cdk473, Cdk474, Cdk475, Cdk476, Cdk477, Cdk478, Cdk479, Cdk480, Cdk481, Cdk482, Cdk483, Cdk484, Cdk485, Cdk486, Cdk487, Cdk488, Cdk489, Cdk490, Cdk491, Cdk492, Cdk493, Cdk494, Cdk495, Cdk496, Cdk497, Cdk498, Cdk499, Cdk500, Cdk501, Cdk502, Cdk503, Cdk504, Cdk505, Cdk506, Cdk507, Cdk508, Cdk509, Cdk510, Cdk511, Cdk512, Cdk513, Cdk514, Cdk515, Cdk516, Cdk517, Cdk518, Cdk519, Cdk520, Cdk521, Cdk522, Cdk523, Cdk524, Cdk525, Cdk526, Cdk527, Cdk528, Cdk529, Cdk530, Cdk531, Cdk532, Cdk533, Cdk534, Cdk535, Cdk536, Cdk537, Cdk538, Cdk539, Cdk540, Cdk541, Cdk542, Cdk543, Cdk544, Cdk545, Cdk546, Cdk547, Cdk548, Cdk549, Cdk550, Cdk551, Cdk552, Cdk553, Cdk554, Cdk555, Cdk556, Cdk557, Cdk558, Cdk559, Cdk560, Cdk561, Cdk562, Cdk563, Cdk564, Cdk565, Cdk566, Cdk567, Cdk568, Cdk569, Cdk570, Cdk571, Cdk572, Cdk573, Cdk574, Cdk575, Cdk576, Cdk577, Cdk578, Cdk579, Cdk580, Cdk581, Cdk582, Cdk583, Cdk584, Cdk585, Cdk586, Cdk587, Cdk588, Cdk589, Cdk590, Cdk591, Cdk592, Cdk593, Cdk594, Cdk595, Cdk596, Cdk597, Cdk598, Cdk599, Cdk600, Cdk601, Cdk602, Cdk603, Cdk604, Cdk605, Cdk606, Cdk607, Cdk608, Cdk609, Cdk610, Cdk611, Cdk612, Cdk613, Cdk614, Cdk615, Cdk616, Cdk617, Cdk618, Cdk619, Cdk620, Cdk621, Cdk622, Cdk623, Cdk624, Cdk625, Cdk626, Cdk627, Cdk628, Cdk629, Cdk630, Cdk631, Cdk632, Cdk633, Cdk634, Cdk635, Cdk636, Cdk637, Cdk638, Cdk639, Cdk640, Cdk641, Cdk642, Cdk643, Cdk644, Cdk645, Cdk646, Cdk647, Cdk648, Cdk649

```
pathview(gene.data=foldchanges, pathway.id="hsa03030")
```

13

Info: Working in directory /Users/chelseazhong/Desktop/bimm143/class14

Info: Writing image file hsa03030.pathview.png



Focus on top 5 upregulated pathways here for demo purposes only

```
keggrespathways <- rownames(keggres$greater)[1:5]
```

```
# Extract the 8 character long IDs part of each string
keggresids = substr(keggrespathways, start=1, stop=8)
keggresids
```

```
[1] "hsa04640" "hsa04630" "hsa00140" "hsa04142" "hsa04330"
```

```
pathview(gene.data=foldchanges, pathway.id=keggresids, species="hsa")
```

```
'select()' returned 1:1 mapping between keys and columns
```

```
Info: Working in directory /Users/chelseazhong/Desktop/bimm143/class14
```

```
Info: Writing image file hsa04640.pathview.png
```

```
'select()' returned 1:1 mapping between keys and columns
```

```
Info: Working in directory /Users/chelseazhong/Desktop/bimm143/class14
```

```
Info: Writing image file hsa04630.pathview.png
```

```
'select()' returned 1:1 mapping between keys and columns
```

```
Info: Working in directory /Users/chelseazhong/Desktop/bimm143/class14
```

```
Info: Writing image file hsa00140.pathview.png
```

```
'select()' returned 1:1 mapping between keys and columns
```

```
Info: Working in directory /Users/chelseazhong/Desktop/bimm143/class14
```

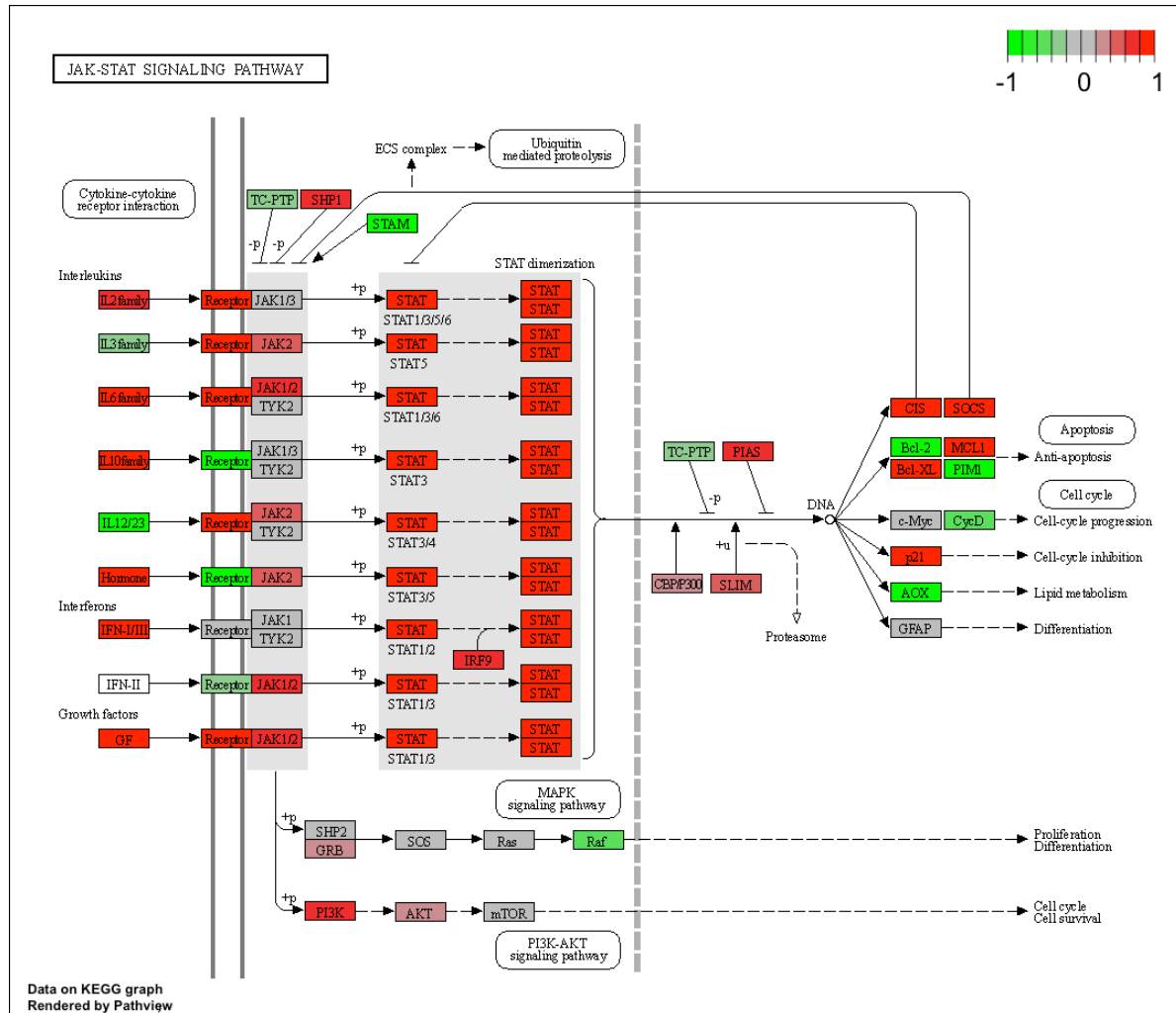
```
Info: Writing image file hsa04142.pathview.png
```

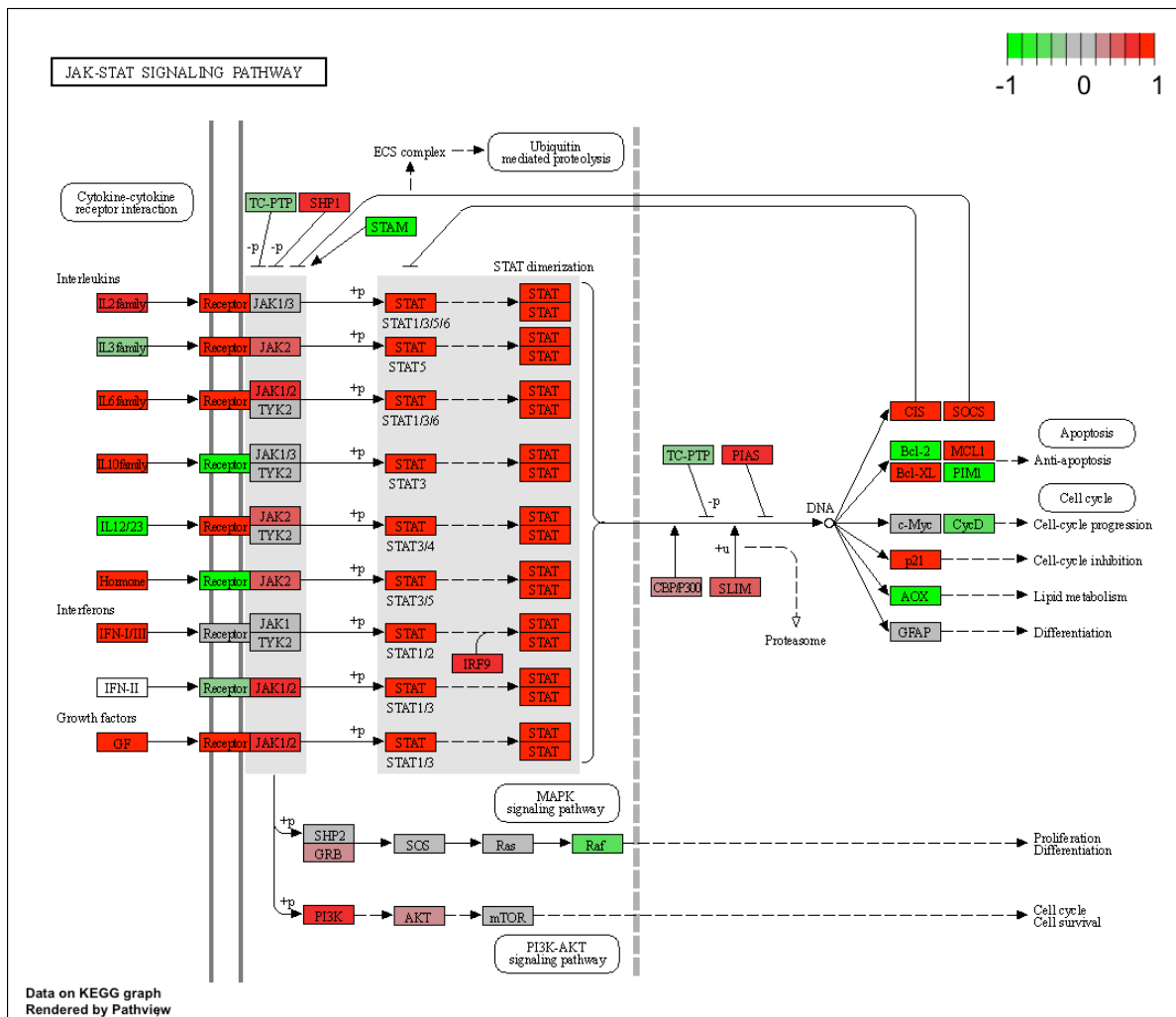
```
Info: some node width is different from others, and hence adjusted!
```

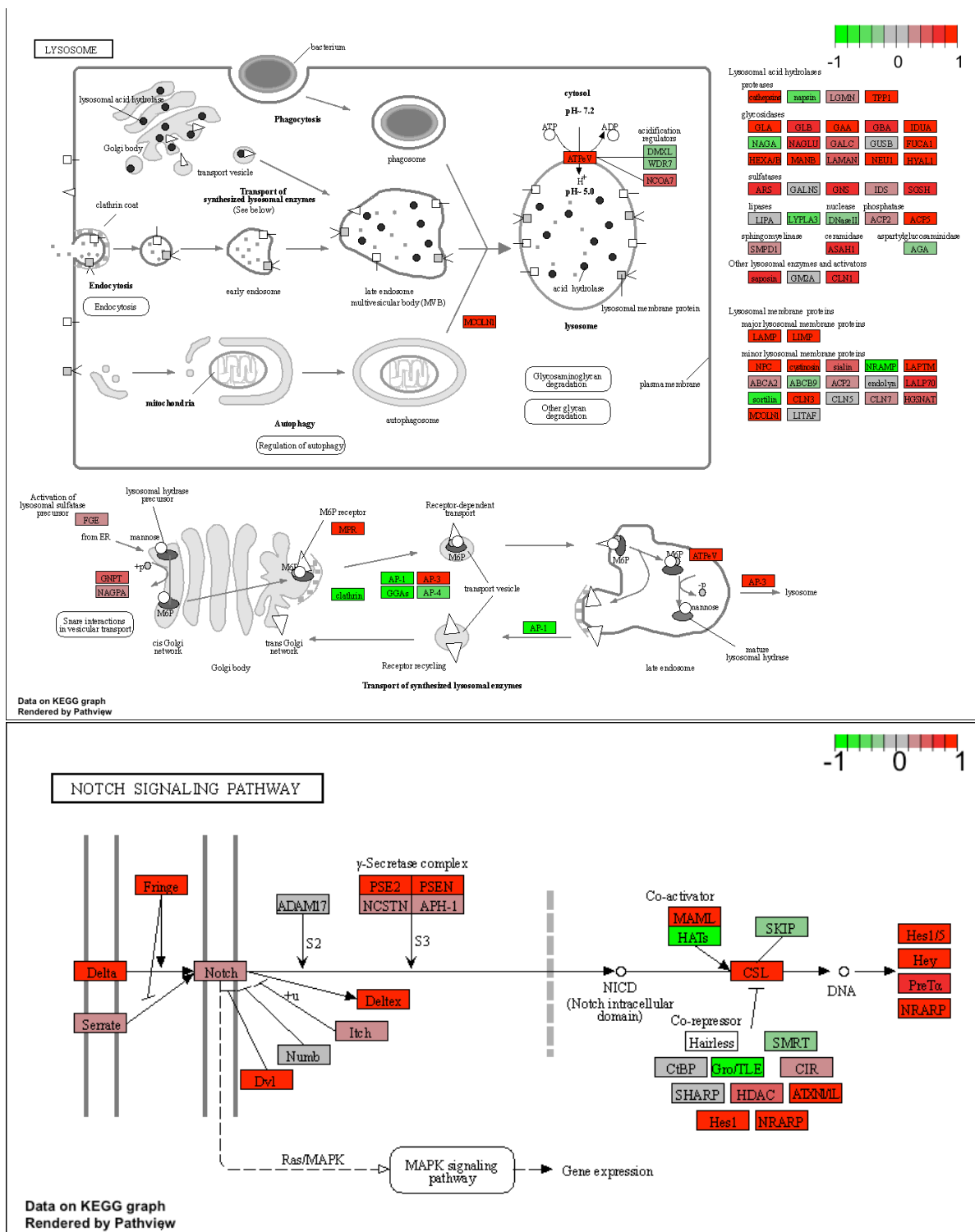
```
'select()' returned 1:1 mapping between keys and columns
```

Info: Working in directory /Users/chelseazhong/Desktop/bimm143/class14

Info: Writing image file hsa04330.pathview.png







gene oncology we can do more analysis

```

data(go.sets.hs)
data(go.subs.hs)

# Focus on Biological Process subset of GO
gobpsets = go.sets.hs[go.subs.hs$BP]

gobpres = gage(foldchanges, gsets=gobpsets, same.dir=TRUE)

lapply(gobpres, head)

```

\$greater

	p.geomean	stat.mean	p.val
GO:0007156 homophilic cell adhesion	8.519724e-05	3.824205	8.519724e-05
GO:0002009 morphogenesis of an epithelium	1.396681e-04	3.653886	1.396681e-04
GO:0048729 tissue morphogenesis	1.432451e-04	3.643242	1.432451e-04
GO:0007610 behavior	1.925222e-04	3.565432	1.925222e-04
GO:0060562 epithelial tube morphogenesis	5.932837e-04	3.261376	5.932837e-04
GO:0035295 tube development	5.953254e-04	3.253665	5.953254e-04

	q.val	set.size	expl
GO:0007156 homophilic cell adhesion	0.1952430	113	8.519724e-05
GO:0002009 morphogenesis of an epithelium	0.1952430	339	1.396681e-04
GO:0048729 tissue morphogenesis	0.1952430	424	1.432451e-04
GO:0007610 behavior	0.1968058	426	1.925222e-04
GO:0060562 epithelial tube morphogenesis	0.3566193	257	5.932837e-04
GO:0035295 tube development	0.3566193	391	5.953254e-04

\$less

	p.geomean	stat.mean	p.val
GO:0048285 organelle fission	1.536227e-15	-8.063910	1.536227e-15
GO:0000280 nuclear division	4.286961e-15	-7.939217	4.286961e-15
GO:0007067 mitosis	4.286961e-15	-7.939217	4.286961e-15
GO:0000087 M phase of mitotic cell cycle	1.169934e-14	-7.797496	1.169934e-14
GO:0007059 chromosome segregation	2.028624e-11	-6.878340	2.028624e-11
GO:0000236 mitotic prometaphase	1.729553e-10	-6.695966	1.729553e-10

	q.val	set.size	expl
GO:0048285 organelle fission	5.843127e-12	376	1.536227e-15
GO:0000280 nuclear division	5.843127e-12	352	4.286961e-15
GO:0007067 mitosis	5.843127e-12	352	4.286961e-15
GO:0000087 M phase of mitotic cell cycle	1.195965e-11	362	1.169934e-14
GO:0007059 chromosome segregation	1.659009e-08	142	2.028624e-11
GO:0000236 mitotic prometaphase	1.178690e-07	84	1.729553e-10

\$stats

	stat.mean	exp1
G0:0007156 homophilic cell adhesion	3.824205	3.824205
G0:0002009 morphogenesis of an epithelium	3.653886	3.653886
G0:0048729 tissue morphogenesis	3.643242	3.643242
G0:0007610 behavior	3.565432	3.565432
G0:0060562 epithelial tube morphogenesis	3.261376	3.261376
G0:0035295 tube development	3.253665	3.253665

```
sig_genes <- res[res$padj <= 0.05 & !is.na(res$padj), "symbol"]  
print(paste("Total number of significant genes:", length(sig_genes)))
```

```
[1] "Total number of significant genes: 8147"
```

```
write.table(sig_genes, file="significant_genes.txt", row.names=FALSE, col.names=FALSE, quo
```

#perform pathway analysis online on significant_gene.txt go to the Reactome website (<https://reactome.org/PathwayBrowser/#TOOL=AT>). Q: What pathway has the most significant “Entities p-value”? Do the most significant pathways listed match your previous KEGG results? What factors could cause differences between the two methods?

Cell cycle, mitotic has the most significant “Entities p-value”, which matches with my previous KEGG results.