# Class 13: RNASeq mini project

Henry (A16354124)

**Background**

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         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
```

```
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
```

```
# Filter count data where you have 0 read count across all samples.
countData = countData[rowSums(countData) > 0, ]
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
```

**Setup for DESeq**

```
library("DESeq2")
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)
```

**Run DESeq**

```
dds <- DESeq(dds)
```

using pre-existing size factors

estimating dispersions

found already estimated dispersions, replacing these

gene-wise dispersion estimates

mean-dispersion relationship

final dispersion estimates

fitting model and testing

**Get results**

```
res <- results(dds)
```

```
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
```

**Add annotation**

```
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  0.3248215    0.551863 5.81042e-01
ENSG00000187634  183.229650      0.4264571  0.1402658    3.040350 2.36304e-03
ENSG00000188976 1651.188076     -0.6927205  0.0548465  -12.630156 1.43993e-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  0.5215598    1.040744 2.97994e-01
ENSG00000188290  108.922128      2.0570638  0.1969053   10.446970 1.51281e-25
ENSG00000187608  350.716868      0.2573837  0.1027266    2.505522 1.22271e-02
ENSG00000188157 9128.439422      0.3899088  0.0467164    8.346302 7.04333e-17
ENSG00000237330    0.158192      0.7859552  4.0804729    0.192614 8.47261e-01
                      padj      symbol      entrez                        name
                 <numeric> <character> <character>                 <character>
ENSG00000279457 6.86555e-01          NA          NA                          NA
ENSG00000187634 5.15718e-03      SAMD11      148398 sterile alpha motif ..
ENSG00000188976 1.76553e-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 4.21970e-16       AGRN     375790                   agrin
ENSG00000237330          NA      RNF223     401934 ring finger protein ..
```

**Visualize results**

```
library(ggplot2)

mycols <- rep("gray", nrow(res) )
mycols[ abs(res$log2FoldChange) > 2 ] <- "red"
inds <- which(res$pvalue < 0.05 & abs(res$log2FoldChange) > 2)
mycols[inds] <- "blue"

ggplot(res) +
  aes(x=log2FoldChange, y=-log10(pvalue)) +
  geom_point(alpha=0.4, color = mycols) +
  theme_minimal() +
  xlab("Log2 Fold Change") +
  ylab("-Log10 p-value") +
  geom_vline(xintercept=c(-2, 2), col="red", linetype="dashed") +
  geom_hline(yintercept=-log10(0.05), col="red", linetype="dashed") +
  ggtitle("Volcano Plot of Differential Expression")
```

## Volcano Plot of Differential Expression



**Pathway analysis**

```r
library(gage)
```

```r
library(gageData)
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).
##############################################################################
```

```
data(kegg.sets.hs)
data(sigmet.idx.hs)
kegg.sets.hs = kegg.sets.hs[sigmet.idx.hs]
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)
```

```
      <NA>       148398        26155       339451        84069        84808
  0.17925708   0.42645712  -0.69272046   0.72975561   0.04057653   0.54281049
```

```
keggres = gage(foldchanges, gsets=kegg.sets.hs)
pathview(gene.data=foldchanges, pathway.id="hsa04110")
```

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

Info: Working in directory /Users/hix/Documents/BGGN213/class13

Info: Writing image file hsa04110.pathview.png

Another version

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

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

Warning: reconcile groups sharing member nodes!

```
      [,1] [,2]
[1,] "9"  "300"
[2,] "9"  "306"
```

Info: Working in directory /Users/hix/Documents/BGGN213/class13

Info: Writing image file hsa04110.pathview.pdf

```
## 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/hix/Documents/BGGN213/class13

Info: Writing image file hsa04640.pathview.png

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

Info: Working in directory /Users/hix/Documents/BGGN213/class13

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

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

Info: Working in directory /Users/hix/Documents/BGGN213/class13

Info: Writing image file hsa00140.pathview.png

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

Info: Working in directory /Users/hix/Documents/BGGN213/class13

Info: Writing image file hsa04142.pathview.png

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

Info: Working in directory /Users/hix/Documents/BGGN213/class13

Info: Writing image file hsa04330.pathview.png
```
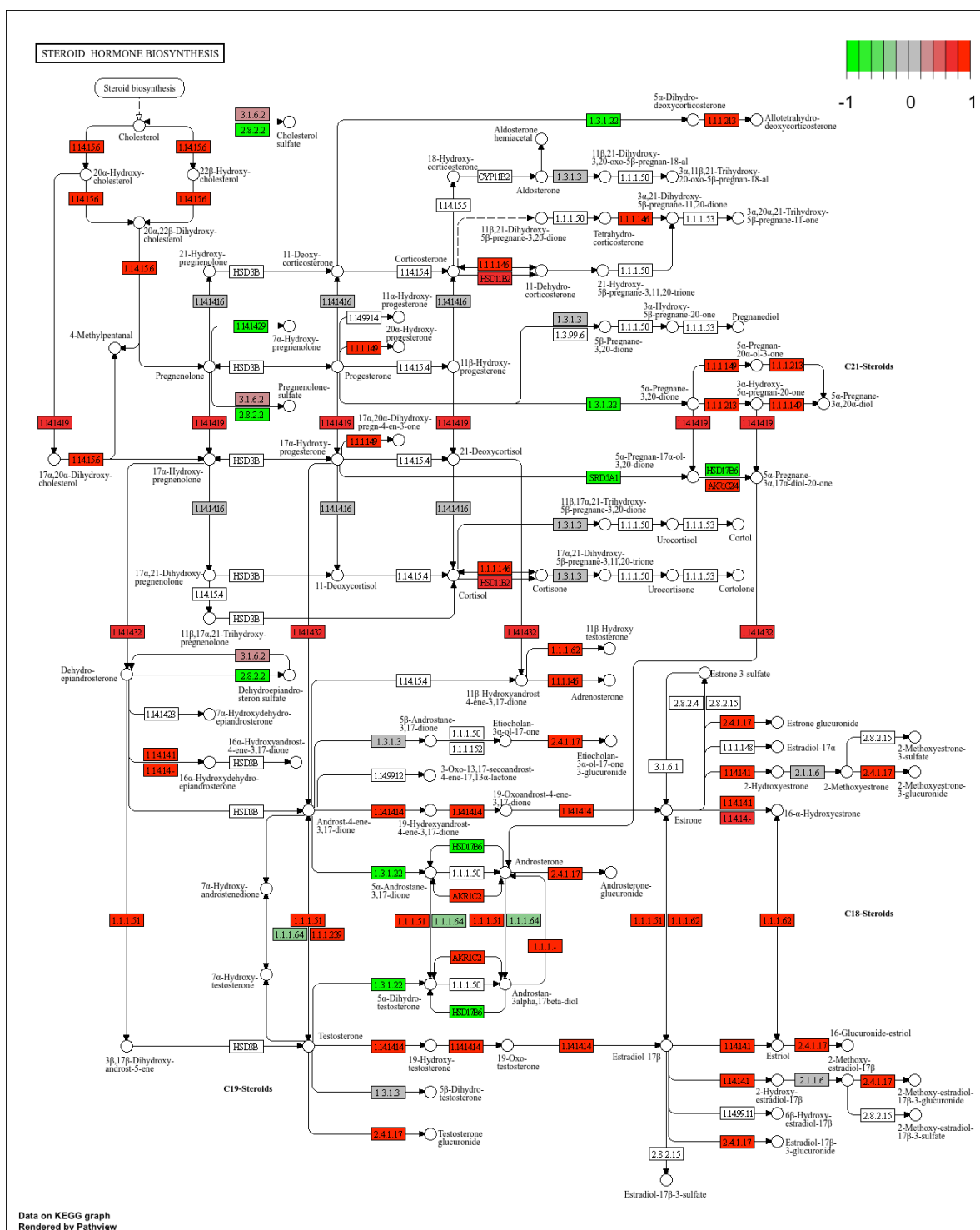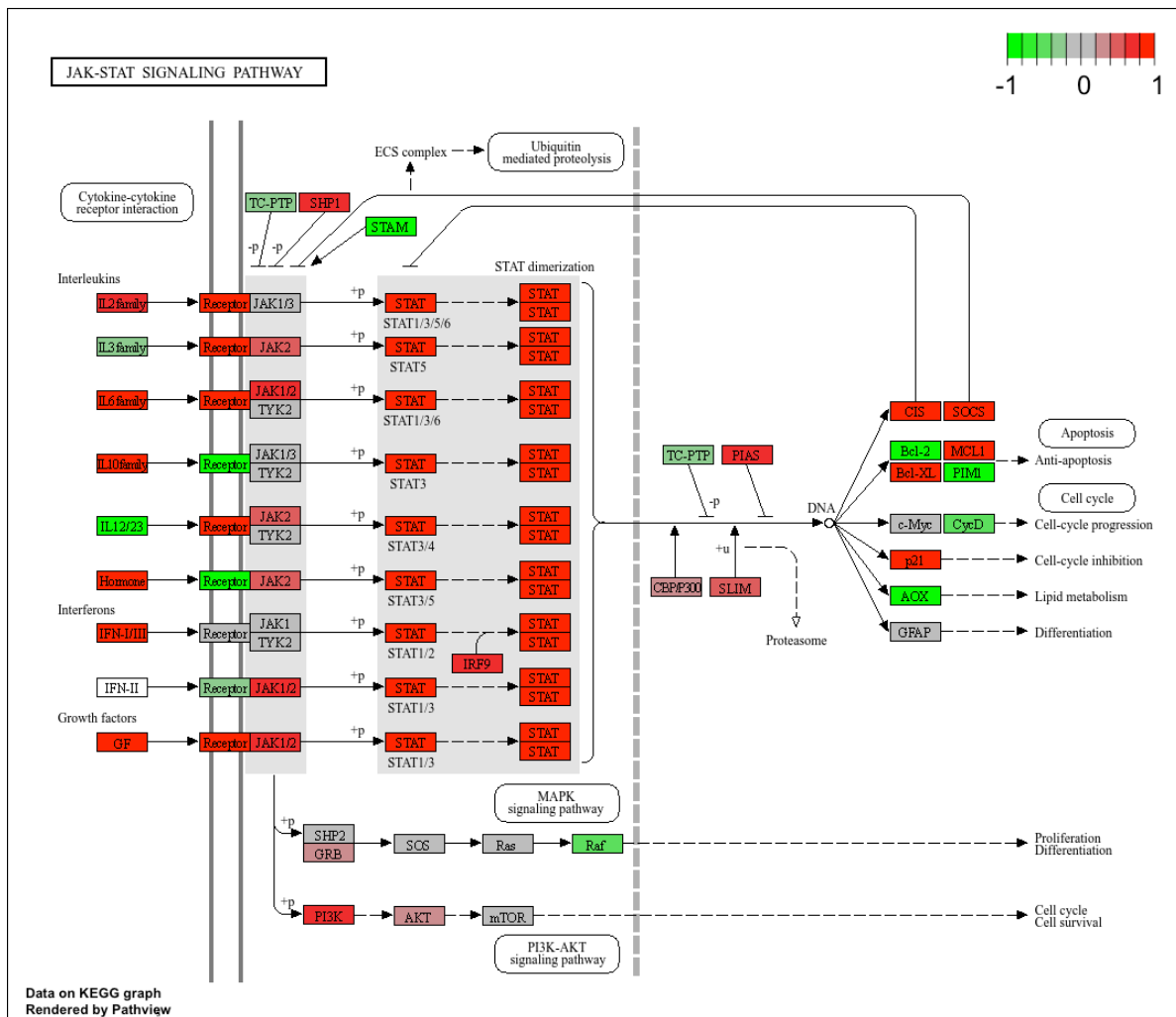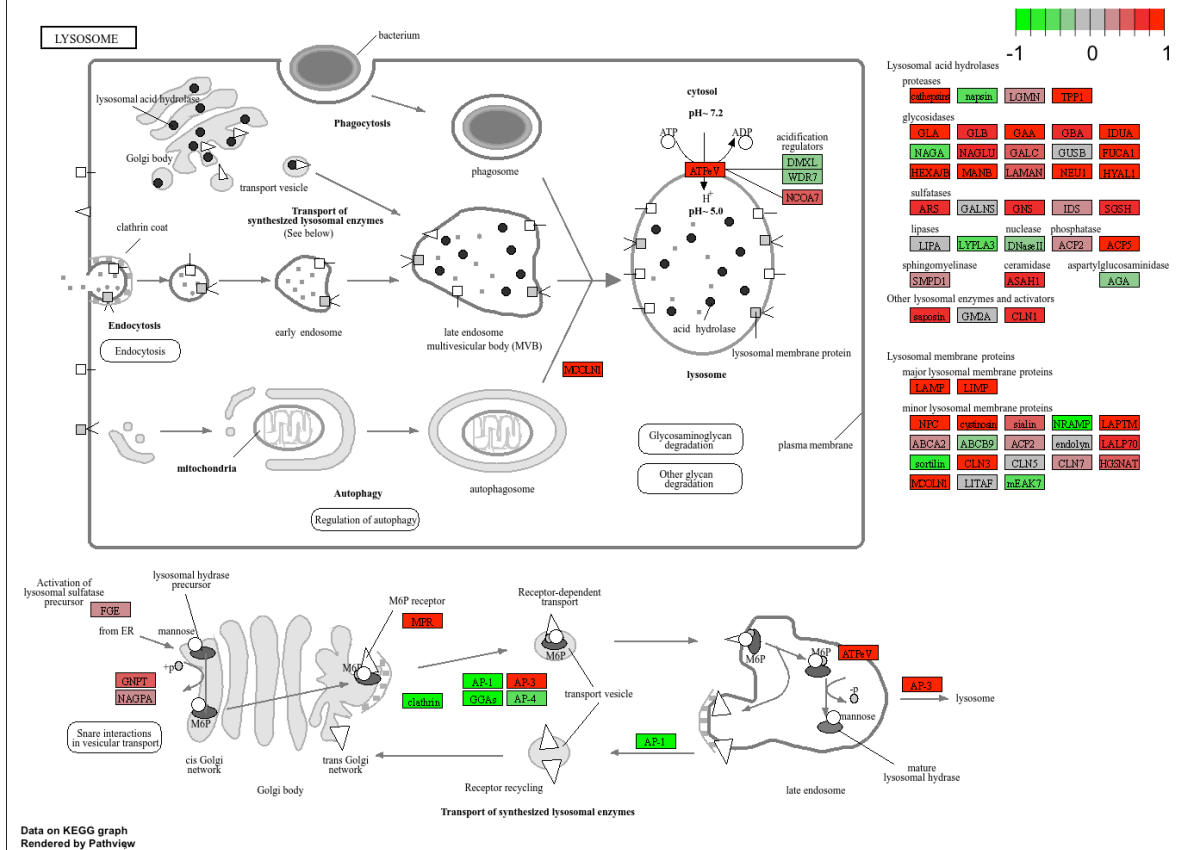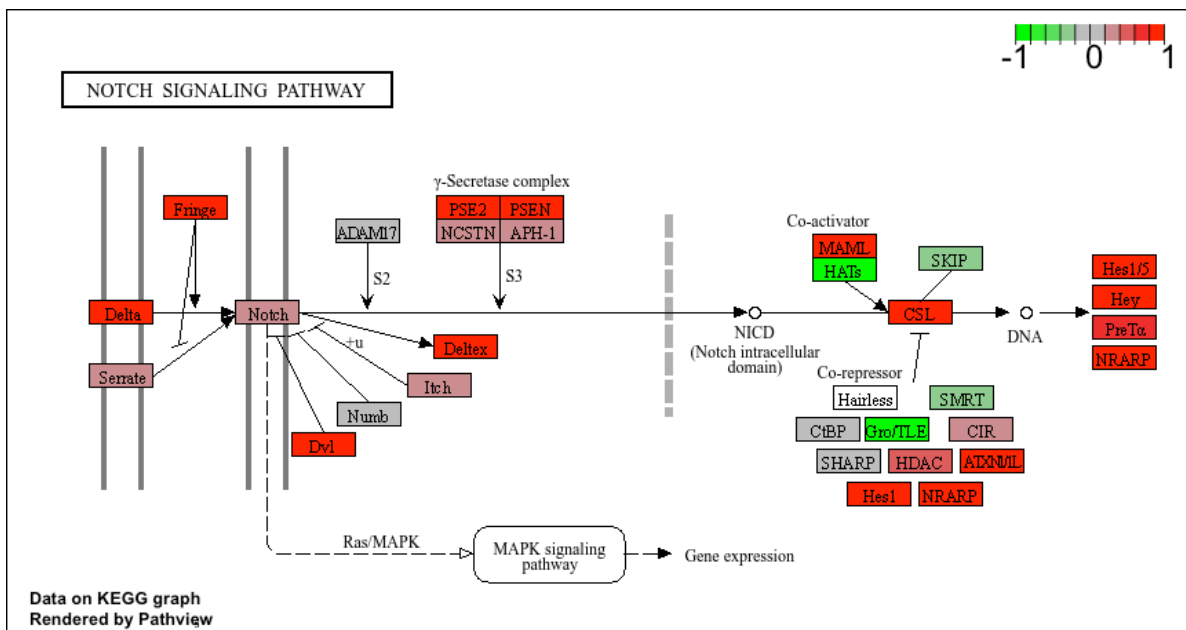
Lymphoid Related Dendritic cell

1    0    1

Thymus

γδ T cell

IL-7

SCF
IL-7

SCF
IL-7

(IL-7)

CD8 T cell

Pro T cell
(DN2)

DN3

DN4

Intermediate
single-positive
cell (ISP)

Double-positive
cell (DP)

CD4 T cell

(CD2)   (CD5)
CD7     CD25
CD38    CD44
(CD71)  CD117
CD127   TdT
HLA-DR

CD2     CD5
CD7     CD25
CD38    CD44
CD71    CD117
(CD127) TdT

CD1     CD2
CD4     CD5
CD7     CD38
(CD44)  (CD117)
TdT

CD2
CD4or8
CD5
CD7
CD38

CD2     CD3
CD4or8  CD5
CD7

CD3
CD5

Regulatory T cell

NKT cell

| SCF | IL-7 |
| --- | --- |

| HLA-DR | CD44 | CD117 | CD25 | CD127 | TdT | CD71 | CD38 | CD7 | CD2 | CD5 | CD1 | CD4 | CD8 | CD3 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |

SCF
IL-7

NK cell Precursor

SCF
IL-7

NK cell

Lymphoid
stem cell,
Double-negative
cell (DN1)

IL-7

CD34
CD44
CD117
TdT
HLA-DR

Pro B Cell

Pre B I cell

Pre B II cell

Immature B cell

B Cell

(CD9)   (CD10)
CD19    (CD20)
CD22    CD24
CD38    CD117
CD127   TdT
TdT     HLA-DR

CD9     CD10
CD19    CD20
CD22    CD24
CD38    CD117
CD127   TdT
HLA-DR

(CD9)   CD19
CD20    CD21
CD22    CD24
CD37    HLA-DR
IgM

(CD5)   (CD9)
CD19    CD19
CD20    CD21
CD22    CD24
(CD23)  CD37
CD35    IgM
HLA-DR
IgD

| IL-7 |
| --- |

| TdT | CD117 | CD10 | CD38 | CD127 | CD9 | HLA-DR | CD19 | CD22 | CD24 | CD25 | CD20 | CD21 | CD37 | IgM | CD23 | CD35 | IgD |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |

Hematopoietic
stem cell

CD34
CD135

SCF
IL-7

| SCF | IL-7 |
| --- | --- |

| CD34 | CD135 | TdT | HLA-DR |
| --- | --- | --- | --- |

SCF
IL-3
IL-4

SCF
IL-4

CFU-Mast

Mast cell

| SCF | IL-3 | IL-4 |
| --- | --- | --- |

SCF
GM-CSF  IL-3

GM-CSF
IL-3

GM-CSF
IL-3

GM-CSF
IL-3

CFU-Bas

Myeloblast

Basophilic
Myelocyte

Basophil

| SCF | IL-3 | GM-CSF |
| --- | --- | --- |

Flt3L   GM-CSF
SCF     IL-3

GM-CSF
IL-3
IL-5

GM-CSF
IL-3
IL-5

GM-CSF
IL-5

CFU-E0

Myeloblast

Eosinophilic
Myelocyte

Eosinophil

| Flt3L | SCF | IL-3 | GM-CSF | IL-5 |
| --- | --- | --- | --- | --- |

Flt3L
CSF
GM-CSF

IL-3
TNF

Flt3L   GM-CSF
SCF     IL-4    TNF

GM-CSF
IL-4

Myeloid Related
Dendritic Cell

CFU-M/DC

GM-CSF
M-CSF
IL-3

GM-CSF
M-CSF
IL-3

GM-CSF
M-CSF
IL-3

GM-CSF
IL-4

Monoblast

Promonocyte

Monocyte

GM-CSF
M-CSF

Macrophage

CD11b   CD13
CD14    CD15
CD33    CD64
CD115   CD116
CD123   CD124
CD126   HLA-DR

CD11b   CD13
CD14    CD33
CD64    CD115
CD116   CD123
CD124   CD126
HLA-DR

CD11b
CD14
CD33
CD64

| Flt3L | SCF | IL-3 | GM-CSF | TNF | IL-4 | M-SCF |
| --- | --- | --- | --- | --- | --- | --- |

| HLA-DR | CD116 | CD123 | CD33 | CD124 | CD126 | CD64 | CD115 | CD13 | CD11b | CD14 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |

Flt3L
SCF
G-CSF
IL-1
IL-3
IL-6
IL-11

Flt3L
SCF
GM-CSF
G-CSF
IL-3

GM-CSF
G-CSF
IL-3

Flt3L
SCF
GM-CSF
IL-3

G-CSF
IL-3

GM-CSF
G-CSF

GM-CSF
G-CSF

GM-CSF
G-CSF

Myeloid
Stem Cell

CFU-GEMM

CFU-GM

CFU-G

Myeloblast

Neutrophilic
Myelocyte

Neutrophil

Bone marrow

CD33    CD34
CD116   CD114
CD121   CD123
CD126   EPOR
IL-9R
HLA-DR

CD15    CD33
CD34    CD64
CD114   CD115
CD116   CD121
CD123   CD124
CD125   CD126
HLA-DR

CD13    CD15
CD33    CD114
CD116   CD121
CD123   CD124
CD125   CD126
HLA-DR

CD13    CD15
CD33    CD114
CD116   CD123
CD124
CD125   CD126

CD11b   CD15
CD33    CD123
CD124   CD125

CD11b
CD15
CD33

| Flt3L | SCF | G-SCF | IL-3 | IL-6 | IL-11 | IL-1 | GM-CSF |
| --- | --- | --- | --- | --- | --- | --- | --- |

| Flt3L | SCF | IL-3 | GM-CSF | G-CSF |
| --- | --- | --- | --- | --- |

| IL-9R | CD34 | HLA-DR | CD116 | CD121 | CD114 | CD123 | CD124 | CD126 | CD33 | CD13 | CD125 | CD11b |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |

Flt3L
SCF
GM-CSF

IL-3
IL-4

SCF
GM-CSF

IL-3
IL-4    EPO

TPO
EPO

EPO

BFU-E

CFU-E

Proerythroblast

Erythrocyte

CD33    CD34
CD117   CD123
EPOR    HLA-DR

CD36
CD235a

CD235a

CD35    CD44
CD55    CD59
CD235a

| Flt3L | SCF | GM-CSF | IL-3 | IL-4 | EPO | TPO |
| --- | --- | --- | --- | --- | --- | --- |

| HLA-DR | EPOR | CD33 | CD34 | CD117 | CD123 | CD36 | CD235a | CD35 | CD44 | CD55 | CD59 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |

Flt3L
SCF
GM-CSF
IL-3

IL-6
IL-11
TPO

Flt3L   Meg-CSF
SCF     IL-3
IL-6    IL-11
TPO

SCF
GM-CSF
IL-3

IL-6
IL-11
TPO

IL-6
IL-11
TPO

BFU-MK

CFU-MK

Mega-
karyocyte

Platelets

CD33    CD34
CD116   CD123
CD126   IL-11R
HLA-DR

CD61
CD116
CD122
CD126

CD9     CD14
CD36    CD41
CD42    CD61
CD61    CD123
CD116
CD126

CD9     CD14
CD36    CD41
CD42    CD49
CD61    CD126

| Flt3L | SCF | IL-3 | IL-6 | IL-11 | GM-CSF | Meg-CSF | TPO |
| --- | --- | --- | --- | --- | --- | --- | --- |

| HLA-DR | CD33 | CD34 | IL-11R | CD116 | CD123 | CD126 | CD61 | CD9 | CD14 | CD36 | CD41 | CD42 | CD49 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |

Data on KEGG graph
Rendered by Pathview

12

STEROID HORMONE BIOSYNTHESIS

Steroid biosynthesis

Data on KEGG graph
Rendered by Pathview

JAK-STAT SIGNALING PATHWAY

-1    0    1

ECS complex → Ubiquitin mediated proteolysis

Cytokine-cytokine receptor interaction

TC-PTP   SHP1
STAM
-p   -p

STAT dimerization

Interleukins

IL2family → Receptor JAK1/3 — +p → STAT   STAT STAT
STAT1/3/5/6

IL3family → Receptor JAK2 — +p → STAT   STAT STAT
STAT5

IL6family → Receptor JAK1/2 TYK2 — +p → STAT   STAT STAT
STAT1/3/6

IL10family → Receptor JAK1/3 TYK2 — +p → STAT   STAT STAT
STAT3

IL12/23 → Receptor JAK2 TYK2 — +p → STAT   STAT STAT
STAT3/4

Hormone → Receptor JAK2 — +p → STAT   STAT STAT
STAT3/5

Interferons

IFN-I/III → Receptor JAK1 TYK2 — +p → STAT   STAT STAT
STAT1/2        IRF9

IFN-II → Receptor JAK1/2 — +p → STAT   STAT STAT
STAT1/3

Growth factors

GF → Receptor JAK1/2 — +p → STAT   STAT STAT
STAT1/3

DNA

TC-PTP   PIAS
-p
CBP/P300   SLIM
+u

Proteasome

CIS   SOCS

Bcl-2   MCL1        → Apoptosis
Bcl-XL  PIM1        Anti-apoptosis

c-Myc   CycD        → Cell cycle
                    Cell-cycle progression

p21                 --→ Cell-cycle inhibition

AOX                 --→ Lipid metabolism

GFAP                --→ Differentiation

MAPK signaling pathway

+p → SHP2 GRB → SOS → Ras → Raf --------→ Proliferation
                                          Differentiation

+p → PI3K → AKT → mTOR --------→ Cell cycle
                                Cell survival

PI3K-AKT signaling pathway

Data on KEGG graph
Rendered by Pathview

14

NOTCH SIGNALING PATHWAY

Data on KEGG graph
Rendered by Pathview



LYSOSOME

Data on KEGG graph
Rendered by Pathview

GO Analyses

```
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         exp1
GO:0007156 homophilic cell adhesion        0.1951953      113 8.519724e-05
GO:0002009 morphogenesis of an epithelium  0.1951953      339 1.396681e-04
GO:0048729 tissue morphogenesis           0.1951953      424 1.432451e-04
GO:0007610 behavior                        0.1967577      426 1.925222e-04
GO:0060562 epithelial tube morphogenesis   0.3565320      257 5.932837e-04
GO:0035295 tube development                0.3565320      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         exp1
GO:0048285 organelle fission              5.841698e-12      376 1.536227e-15
GO:0000280 nuclear division               5.841698e-12      352 4.286961e-15
GO:0007067 mitosis                        5.841698e-12      352 4.286961e-15
GO:0000087 M phase of mitotic cell cycle  1.195672e-11      362 1.169934e-14
GO:0007059 chromosome segregation         1.658603e-08      142 2.028624e-11
```

```
GO:0000236 mitotic prometaphase            1.178402e-07        84 1.729553e-10
```

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

```
head(gobpres$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        exp1
GO:0048285 organelle fission          5.841698e-12      376 1.536227e-15
GO:0000280 nuclear division           5.841698e-12      352 4.286961e-15
GO:0007067 mitosis                    5.841698e-12      352 4.286961e-15
GO:0000087 M phase of mitotic cell cycle 1.195672e-11      362 1.169934e-14
GO:0007059 chromosome segregation     1.658603e-08      142 2.028624e-11
GO:0000236 mitotic prometaphase       1.178402e-07       84 1.729553e-10
```
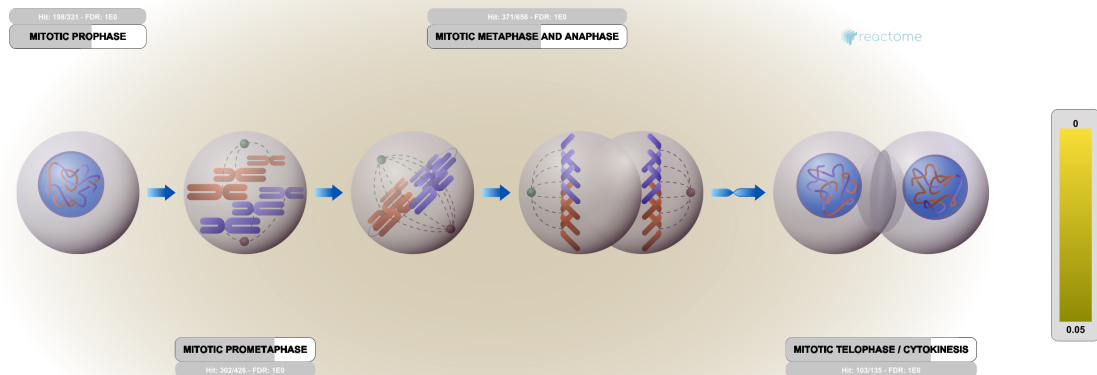
**Reactome**

Some folks like it, we can try it out here

```
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, quote=
```

## Save results

```r
save(res, file = "my_results.RData")
write.csv(res,file = "my_results.csv")
```