

Lab 14: RNASeq

Emily Ignatoff (A16732102)

Today we will analyze differential expression of protein-coding genes in lung fibroblast cells which have lost their **HoxA1** transcription factor (Trapnell et al, 2013)

Differential Expression Analysis

First, let's load in our data:

```
library(DESeq2)
```

Warning: package 'DESeq2' was built under R version 4.3.3

Warning: package 'matrixStats' was built under R version 4.3.3

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

```
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					

We will need to remove `countData$length` in order to have the `colData` and `countData` match.

Q. Complete the code below to remove the troublesome first column from `countData`.

```
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

We can also clean up our data by removing all the genes where there is 0 expression all together.

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).

```
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

We can now set up our dds objects for DESeq and run the DESeq pipeline:

```
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

```

We can now specifically obtain result for the HOXA1 knockdown and siRNA control:

```
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)
```

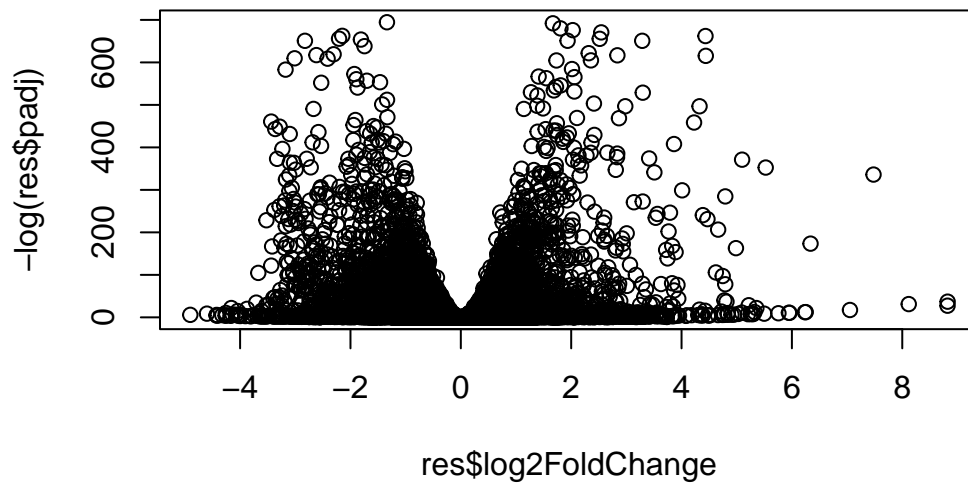
```

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

```

Let's now create a volcano plot to better visualize the up- and down- regulated genes:

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



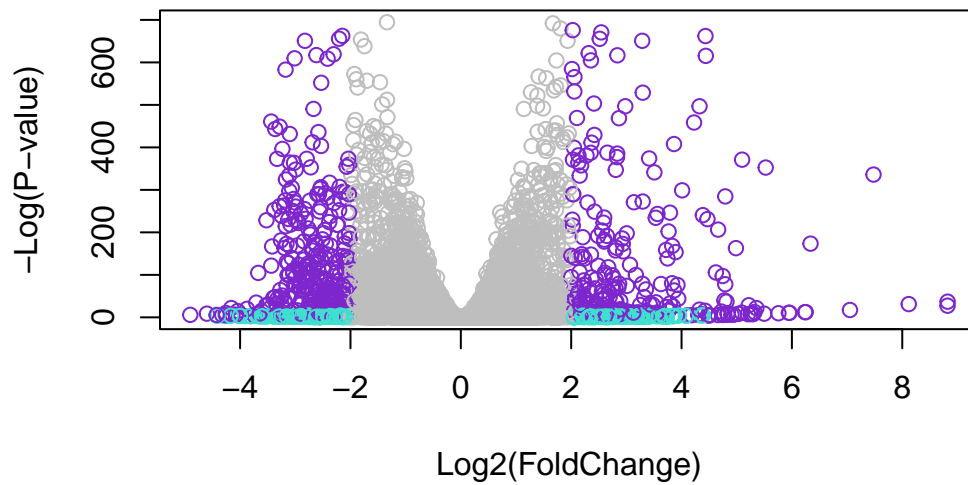
Q. Improve this plot by completing the below code, which adds color and axis labels

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

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

# Color 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 ] <- "purple3"

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



We can also improve this volcano plot with gene annotations, which we can add using the `mapIDs()` function:

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	1.43990e-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	1.51282e-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
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	1.76549e-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	..

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")
```

Pathway Analysis

We will now use the `gage` package to draw pathway diagrams for enriched pathways will elements colored by regulation level and direction.

```
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)
```

Let us focus on signaling and metabolic pathways:

```
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)
```

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

Run the gage pathway analysis:

```
keggres = gage(foldchanges, gsets=kegg.sets.hs)
attributes(keggres)
```

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

Let's look at the first downregulated 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 C:/Users/goose/OneDrive/Documents/Bioinformatics_class/class_14

CELL CYCLE

Cohesion loading

Cohesion establishment

DNA damage checkpoint

Apoptosis

Ungate-mediated proteolysis

MAPK signaling pathway

Growth factor

Growth factor withdrawal

R-point (START)

S-phase proteins, Cyclin

DNA replication

DNA biosynthesis

G1

S

G2

M

Protein interactions:

- Cohesion loading:** ATRX, RFX1, SMC1, SMC3, RFX1, RFX2, RFX3, RFX4, RFX5, RFX6, RFX7, RFX8, RFX9, RFX10, RFX11, RFX12, RFX13, RFX14, RFX15, RFX16, RFX17, RFX18, RFX19, RFX20, RFX21, RFX22, RFX23, RFX24, RFX25, RFX26, RFX27, RFX28, RFX29, RFX30, RFX31, RFX32, RFX33, RFX34, RFX35, RFX36, RFX37, RFX38, RFX39, RFX40, RFX41, RFX42, RFX43, RFX44, RFX45, RFX46, RFX47, RFX48, RFX49, RFX50, RFX51, RFX52, RFX53, RFX54, RFX55, RFX56, RFX57, RFX58, RFX59, RFX60, RFX61, RFX62, RFX63, RFX64, RFX65, RFX66, RFX67, RFX68, RFX69, RFX70, RFX71, RFX72, RFX73, RFX74, RFX75, RFX76, RFX77, RFX78, RFX79, RFX80, RFX81, RFX82, RFX83, RFX84, RFX85, RFX86, RFX87, RFX88, RFX89, RFX90, RFX91, RFX92, RFX93, RFX94, RFX95, RFX96, RFX97, RFX98, RFX99, RFX100, RFX101, RFX102, RFX103, RFX104, RFX105, RFX106, RFX107, RFX108, RFX109, RFX110, RFX111, RFX112, RFX113, RFX114, RFX115, RFX116, RFX117, RFX118, RFX119, RFX120, RFX121, RFX122, RFX123, RFX124, RFX125, RFX126, RFX127, RFX128, RFX129, RFX130, RFX131, RFX132, RFX133, RFX134, RFX135, RFX136, RFX137, RFX138, RFX139, RFX140, RFX141, RFX142, RFX143, RFX144, RFX145, RFX146, RFX147, RFX148, RFX149, RFX150, RFX151, RFX152, RFX153, RFX154, RFX155, RFX156, RFX157, RFX158, RFX159, RFX160, RFX161, RFX162, RFX163, RFX164, RFX165, RFX166, RFX167, RFX168, RFX169, RFX170, RFX171, RFX172, RFX173, RFX174, RFX175, RFX176, RFX177, RFX178, RFX179, RFX180, RFX181, RFX182, RFX183, RFX184, RFX185, RFX186, RFX187, RFX188, RFX189, RFX190, RFX191, RFX192, RFX193, RFX194, RFX195, RFX196, RFX197, RFX198, RFX199, RFX200, RFX201, RFX202, RFX203, RFX204, RFX205, RFX206, RFX207, RFX208, RFX209, RFX210, RFX211, RFX212, RFX213, RFX214, RFX215, RFX216, RFX217, RFX218, RFX219, RFX220, RFX221, RFX222, RFX223, RFX224, RFX225, RFX226, RFX227, RFX228, RFX229, RFX230, RFX231, RFX232, RFX233, RFX234, RFX235, RFX236, RFX237, RFX238, RFX239, RFX240, RFX241, RFX242, RFX243, RFX244, RFX245, RFX246, RFX247, RFX248, RFX249, RFX250, RFX251, RFX252, RFX253, RFX254, RFX255, RFX256, RFX257, RFX258, RFX259, RFX260, RFX261, RFX262, RFX263, RFX264, RFX265, RFX266, RFX267, RFX268, RFX269, RFX270, RFX271, RFX272, RFX273, RFX274, RFX275, RFX276, RFX277, RFX278, RFX279, RFX280, RFX281, RFX282, RFX283, RFX284, RFX285, RFX286, RFX287, RFX288, RFX289, RFX290, RFX291, RFX292, RFX293, RFX294, RFX295, RFX296, RFX297, RFX298, RFX299, RFX300, RFX301, RFX302, RFX303, RFX304, RFX305, RFX306, RFX307, RFX308, RFX309, RFX310, RFX311, RFX312, RFX313, RFX314, RFX315, RFX316, RFX317, RFX318, RFX319, RFX320, RFX321, RFX322, RFX323, RFX324, RFX325, RFX326, RFX327, RFX328, RFX329, RFX330, RFX331, RFX332, RFX333, RFX334, RFX335, RFX336, RFX337, RFX338, RFX339, RFX340, RFX341, RFX342, RFX343, RFX344, RFX345, RFX346, RFX347, RFX348, RFX349, RFX350, RFX351, RFX352, RFX353, RFX354, RFX355, RFX356, RFX357, RFX358, RFX359, RFX360, RFX361, RFX362, RFX363, RFX364, RFX365, RFX366, RFX367, RFX368, RFX369, RFX370, RFX371, RFX372, RFX373, RFX374, RFX375, RFX376, RFX377, RFX378, RFX379, RFX380, RFX381, RFX382, RFX383, RFX384, RFX385, RFX386, RFX387, RFX388, RFX389, RFX390, RFX391, RFX392, RFX393, RFX394, RFX395, RFX396, RFX397, RFX398, RFX399, RFX400, RFX401, RFX402, RFX403, RFX404, RFX405, RFX406, RFX407, RFX408, RFX409, RFX410, RFX411, RFX412, RFX413, RFX414, RFX415, RFX416, RFX417, RFX418, RFX419, RFX420, RFX421, RFX422, RFX423, RFX424, RFX425, RFX426, RFX427, RFX428, RFX429, RFX430, RFX431, RFX432, RFX433, RFX434, RFX435, RFX436, RFX437, RFX438, RFX439, RFX440, RFX441, RFX442, RFX443, RFX444, RFX445, RFX446, RFX447, RFX448, RFX449, RFX450, RFX451, RFX452, RFX453, RFX454, RFX455, RFX456, RFX457, RFX458, RFX459, RFX460, RFX461, RFX462, RFX463, RFX464, RFX465, RFX466, RFX467, RFX468, RFX469, RFX470, RFX471, RFX472, RFX473, RFX474, RFX475, RFX476, RFX477, RFX478, RFX479, RFX480, RFX481, RFX482, RFX483, RFX484, RFX485, RFX486, RFX487, RFX488, RFX489, RFX490, RFX491, RFX492, RFX493, RFX494, RFX495, RFX496, RFX497, RFX498, RFX499, RFX500, RFX501, RFX502, RFX503, RFX504, RFX505, RFX506, RFX507, RFX508, RFX509, RFX510, RFX511, RFX512, RFX513, RFX514, RFX515, RFX516, RFX517, RFX518, RFX519, RFX520, RFX521, RFX522, RFX523, RFX524, RFX525, RFX526, RFX527, RFX528, RFX529, RFX530, RFX531, RFX532, RFX533, RFX534, RFX535, RFX536, RFX537, RFX538, RFX539, RFX540, RFX541, RFX542, RFX543, RFX544, RFX545, RFX546, RFX547, RFX548, RFX549, RFX550, RFX551, RFX552, RFX553, RFX554, RFX555, RFX556, RFX557, RFX558, RFX559, RFX560, RFX561, RFX562, RFX563, RFX564, RFX565, RFX566, RFX567, RFX568, RFX569, RFX570, RFX571, RFX572, RFX573, RFX574, RFX575, RFX576, RFX577, RFX578, RFX579, RFX580, RFX581, RFX582, RFX583, RFX584, RFX585, RFX586, RFX587, RFX588, RFX589, RFX590, RFX591, RFX592, RFX593, RFX594, RFX595, RFX596, RFX597, RFX598, RFX599, RFX600, RFX601, RFX602, RFX603, RFX604, RFX605, RFX606, RFX607, RFX608, RFX609, RFX610, RFX611, RFX612, RFX613, RFX614, RFX615, RFX616, RFX617, RFX618, RFX619, RFX620, RFX621, RFX622, RFX623, RFX624, RFX625, RFX626, RFX627, RFX628, RFX629, RFX630, RFX631, RFX632, RFX633, RFX634, RFX635, RFX636, RFX637, RFX638, RFX639, RFX640, RFX641, RFX642, R

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

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

11

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

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

Info: Working in directory C:/Users/goose/OneDrive/Documents/Bioinformatics_class/class_14

Info: Writing image file hsa04640.pathview.png

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

Info: Working in directory C:/Users/goose/OneDrive/Documents/Bioinformatics_class/class_14

Info: Writing image file hsa04630.pathview.png

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

Info: Working in directory C:/Users/goose/OneDrive/Documents/Bioinformatics_class/class_14

Info: Writing image file hsa00140.pathview.png

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

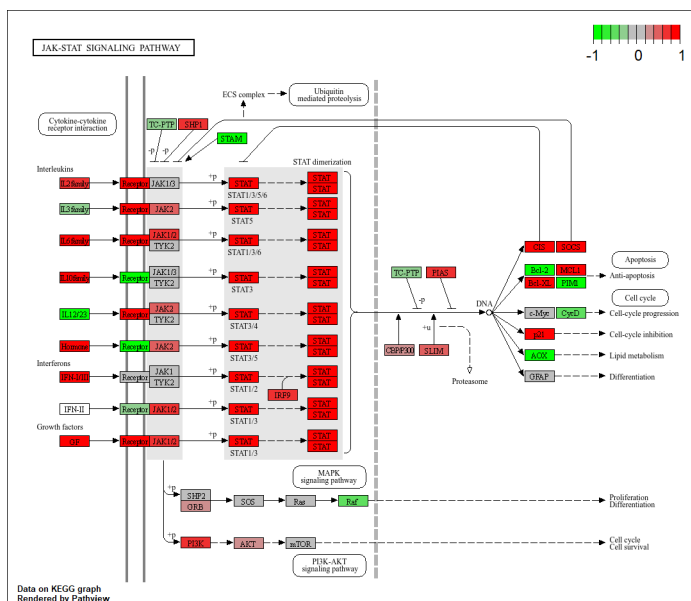
Info: Working in directory C:/Users/goose/OneDrive/Documents/Bioinformatics_class/class_14

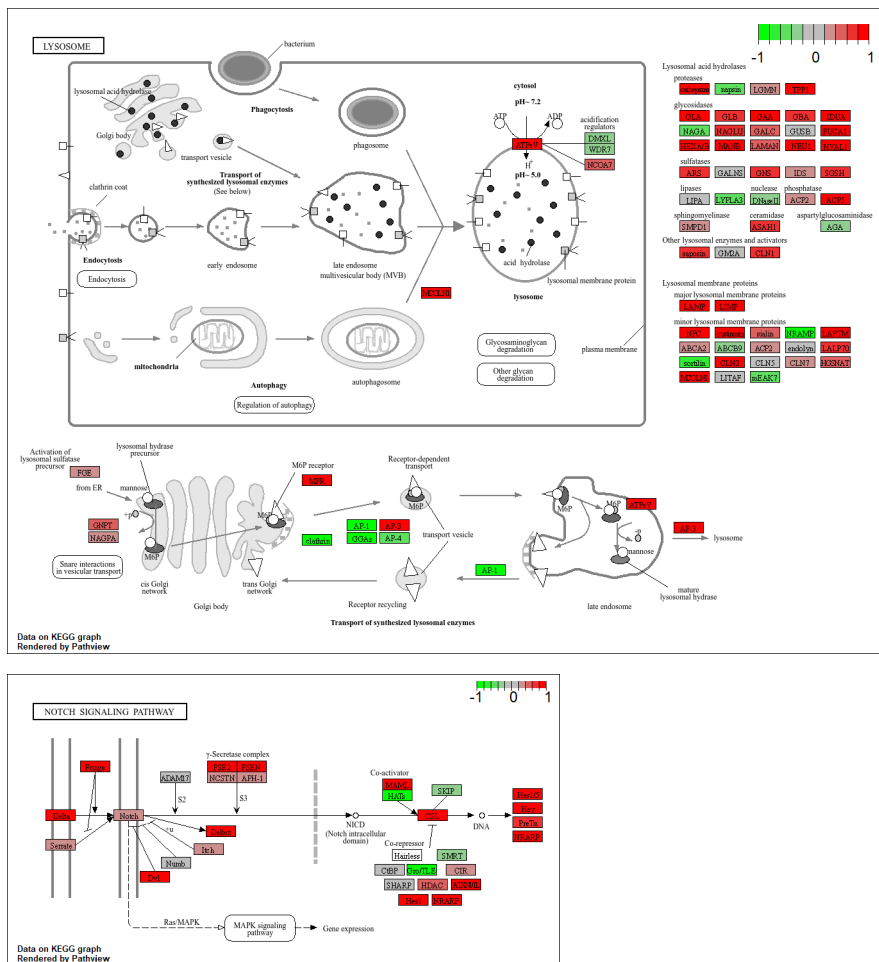
Info: Writing image file hsa04142.pathview.png

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

Info: Working in directory C:/Users/goose/OneDrive/Documents/Bioinformatics_class/class_14

Info: Writing image file hsa04330.pathview.png





Q. Can you do the same procedure as above to plot the pathview figures for the top 5 down-regulated pathways?

```
keggrespathways2 <- rownames(keggres$less)[1:5]
```

```
keggresids2 = substr(keggrespathways2, start=1, stop=8)
keggresids2
```

```
[1] "hsa04110" "hsa03030" "hsa03013" "hsa03440" "hsa04114"
```

```
pathview(gene.data=foldchanges, pathway.id=keggresids2, species="hsa")
```

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

Info: Working in directory C:/Users/goose/OneDrive/Documents/Bioinformatics_class/class_14

Info: Writing image file hsa04110.pathview.png

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

Info: Working in directory C:/Users/goose/OneDrive/Documents/Bioinformatics_class/class_14

Info: Writing image file hsa03030.pathview.png

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

Info: Working in directory C:/Users/goose/OneDrive/Documents/Bioinformatics_class/class_14

Info: Writing image file hsa03013.pathview.png

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

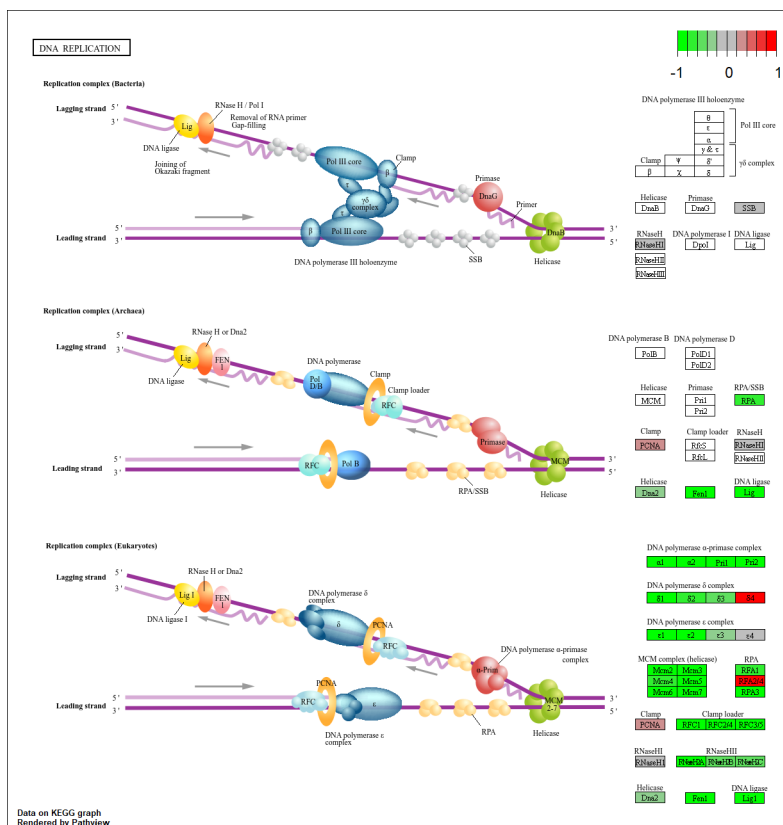
Info: Working in directory C:/Users/goose/OneDrive/Documents/Bioinformatics_class/class_14

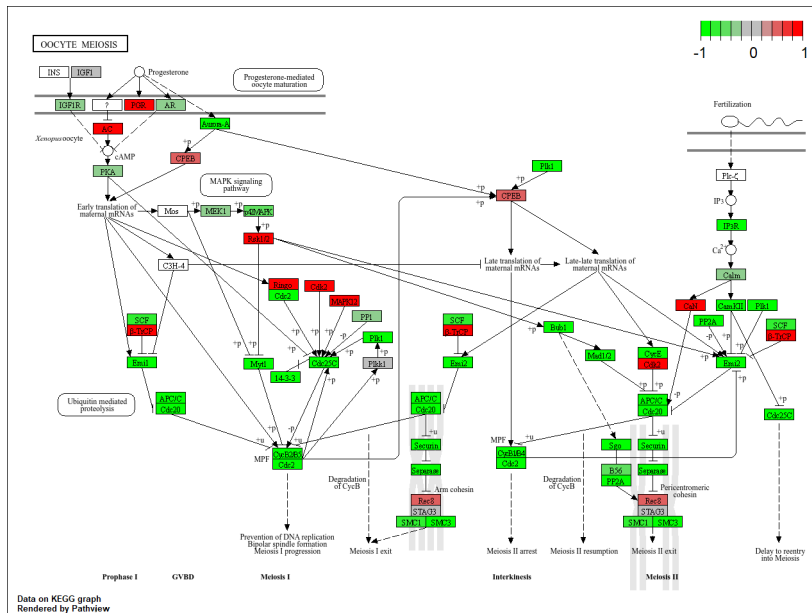
Info: Writing image file hsa03440.pathview.png

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

Info: Working in directory C:/Users/goose/OneDrive/Documents/Bioinformatics_class/class_14

Info: Writing image file hsa04114.pathview.png





Gene Ontology (GO)

We can complete similar analyses to KEGG Pathways with a different dataset, Gene Ontology (GO):

We will be focusing on the **Biological Processes(BP)** section of the gene ontologies

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

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
G0:0007156 homophilic cell adhesion	8.519724e-05	3.824205	8.519724e-05
G0:0002009 morphogenesis of an epithelium	1.396681e-04	3.653886	1.396681e-04
G0:0048729 tissue morphogenesis	1.432451e-04	3.643242	1.432451e-04
G0:0007610 behavior	1.925222e-04	3.565432	1.925222e-04
G0:0060562 epithelial tube morphogenesis	5.932837e-04	3.261376	5.932837e-04
G0:0035295 tube development	5.953254e-04	3.253665	5.953254e-04

	q.val	set.size	exp1
G0:0007156 homophilic cell adhesion	0.1952430	113	8.519724e-05
G0:0002009 morphogenesis of an epithelium	0.1952430	339	1.396681e-04
G0:0048729 tissue morphogenesis	0.1952430	424	1.432451e-04
G0:0007610 behavior	0.1968058	426	1.925222e-04
G0:0060562 epithelial tube morphogenesis	0.3566193	257	5.932837e-04
G0:0035295 tube development	0.3566193	391	5.953254e-04

\$less

	p.geomean	stat.mean	p.val
G0:0048285 organelle fission	1.536227e-15	-8.063910	1.536227e-15
G0:0000280 nuclear division	4.286961e-15	-7.939217	4.286961e-15
G0:0007067 mitosis	4.286961e-15	-7.939217	4.286961e-15
G0:0000087 M phase of mitotic cell cycle	1.169934e-14	-7.797496	1.169934e-14
G0:0007059 chromosome segregation	2.028624e-11	-6.878340	2.028624e-11
G0:0000236 mitotic prometaphase	1.729553e-10	-6.695966	1.729553e-10

	q.val	set.size	exp1
G0:0048285 organelle fission	5.843127e-12	376	1.536227e-15
G0:0000280 nuclear division	5.843127e-12	352	4.286961e-15
G0:0007067 mitosis	5.843127e-12	352	4.286961e-15
G0:0000087 M phase of mitotic cell cycle	1.195965e-11	362	1.169934e-14
G0:0007059 chromosome segregation	1.659009e-08	142	2.028624e-11
G0:0000236 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

Reactome Analysis

The Reactome database consists of biological molecules and their relationships to pathways and processes. Our list of significantly expressed genes as determined above will be used in Reactome for over-representation enrichment analysis and pathway-topology analysis

Generate the significant genes list as a text file:

```
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=
```

We can now upload this `significant_genes.txt` file to Reactome and analyze by projecting these genes onto human genomes

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?

The pathway with the most significant “Entities p-value” is *SLC15A4:TASL-dependent IRF5 activation*, which is not listed in the KEGG results. There are overlapping pathways between KEGG and GO but the Gene Ontology has pathways not present in either of the top five KEGG pathways.

I imagine that the difference in targets between KEGG and GO could lead to these differences. KEGG seeks to understand differential gene expression while GO has several other goals (for example biological processes as we targeted here).