

class12

Working on RNAseq today!

```
library(BiocManager)
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 objects are masked from 'package:base':
```

```
expand.grid, I, unname
```

```
Loading required package: IRanges
```

```
Attaching package: 'IRanges'
```

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

```
windows
```

```
Loading required package: GenomicRanges
```

```
Loading required package: GenomeInfoDb
```

```
Loading required package: SummarizedExperiment
```

```
Loading required package: MatrixGenerics
```

```
Loading required package: matrixStats
```

```
Attaching package: 'MatrixGenerics'
```

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

```
colAlls, colAnyNAs, colAnyNs, colAvgsPerRowSet, 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, rowAnyNs, rowAvgsPerColSet,  
rowCollapse, rowCounts, rowCummaxs, rowCummins, rowCumprods,  
rowCumsums, rowDiffs, rowIQRDiffs, rowIQRs, rowLogSumExps,  
rowMadDiffs, rowMads, rowMaxs, rowMeans2, rowMedians, rowMins,
```

```
rowOrderStats, rowProds, rowQuantiles, rowRanges, rowRanks,  
rowSdDiff, rowSds, rowSums2, rowTabulates, rowVarDiff, 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
```

```
counts <- read.csv("airway_scaledcounts.csv", row.names=1)  
metadata <- read.csv("airway_metadata.csv")
```

```
head(counts)
```

	SRR1039508	SRR1039509	SRR1039512	SRR1039513	SRR1039516
ENSG000000000003	723	486	904	445	1170
ENSG000000000005	0	0	0	0	0
ENSG00000000419	467	523	616	371	582
ENSG00000000457	347	258	364	237	318
ENSG00000000460	96	81	73	66	118
ENSG00000000938	0	0	1	0	2
	SRR1039517	SRR1039520	SRR1039521		
ENSG000000000003	1097	806	604		
ENSG000000000005	0	0	0		

```
ENSG00000000419      781      417      509
ENSG00000000457      447      330      324
ENSG00000000460      94       102      74
ENSG00000000938      0        0        0
```

```
head(metadata)
```

```
  id      dex celltype     geo_id
1 SRR1039508 control   N61311 GSM1275862
2 SRR1039509 treated   N61311 GSM1275863
3 SRR1039512 control   N052611 GSM1275866
4 SRR1039513 treated   N052611 GSM1275867
5 SRR1039516 control   N080611 GSM1275870
6 SRR1039517 treated   N080611 GSM1275871
```

Q1. How many genes are in this dataset? 38694

```
nrow(counts)
```

```
[1] 38694
```

```
ncol(counts)
```

```
[1] 8
```

Let's make sure that the ID column of the metadata match the order of the columns in countData.

```
metadata$id == colnames(counts)
```

```
[1] TRUE TRUE TRUE TRUE TRUE TRUE TRUE TRUE
```

Can see if all the items in the vector are TRUE using `all()`.

```
all(metadata$id == colnames(counts))
```

```
[1] TRUE
```

you can also ask if all all the items are true.

```
#are there any items that are not true?  
!all(metadata$id == colnames(counts))
```

```
[1] FALSE
```

```
#Analysis by hand
```

```
metadata
```

```
      id      dex celltype     geo_id  
1 SRR1039508 control    N61311 GSM1275862  
2 SRR1039509 treated    N61311 GSM1275863  
3 SRR1039512 control    N052611 GSM1275866  
4 SRR1039513 treated    N052611 GSM1275867  
5 SRR1039516 control    N080611 GSM1275870  
6 SRR1039517 treated    N080611 GSM1275871  
7 SRR1039520 control    N061011 GSM1275874  
8 SRR1039521 treated    N061011 GSM1275875
```

Let's first extract our counts for control samples as I want to compare this to the counts for treated (i.e. with drug) samples

```
control inds <- metadata$dex == "control"  
metadata$id[ control inds]
```

```
[1] "SRR1039508" "SRR1039512" "SRR1039516" "SRR1039520"
```

```
control counts <- counts[ control inds]  
head(control counts)
```

	SRR1039508	SRR1039512	SRR1039516	SRR1039520
ENSG000000000003	723	904	1170	806
ENSG000000000005	0	0	0	0
ENSG000000000419	467	616	582	417
ENSG000000000457	347	364	318	330
ENSG000000000460	96	73	118	102
ENSG000000000938	0	1	2	0

Q2. How many ‘control’ cell lines do we have? 4 cell lines

Q3. How would you make the above code in either approach more robust? You would need to use rowMeans to find the average instead of manually dividing the sum by 4 samples.

I want a single summary counts value for each gene in the control experiment. I will start by taking the average.

```
#apply(control.counts,1, mean)
control.mean <-rowMeans(control.counts)
```

Q4. Follow the same procedure for the treated samples (i.e. calculate the mean per gene across drug treated samples and assign to a labeled vector called treated.mean)

```
treated inds <- metadata$dex == "treated"
metadata$id[ treated.inds]
```

```
[1] "SRR1039509" "SRR1039513" "SRR1039517" "SRR1039521"
```

```
treated.counts <- counts[ treated.inds]
head(treated.counts)
```

	SRR1039509	SRR1039513	SRR1039517	SRR1039521
ENSG000000000003	486	445	1097	604
ENSG000000000005	0	0	0	0
ENSG000000000419	523	371	781	509
ENSG000000000457	258	237	447	324
ENSG000000000460	81	66	94	74
ENSG000000000938	0	0	0	0

```
treated.mean <-rowMeans(treated.counts)
```

To help us stay organized, let’s make a new data.frame to stores these results together.

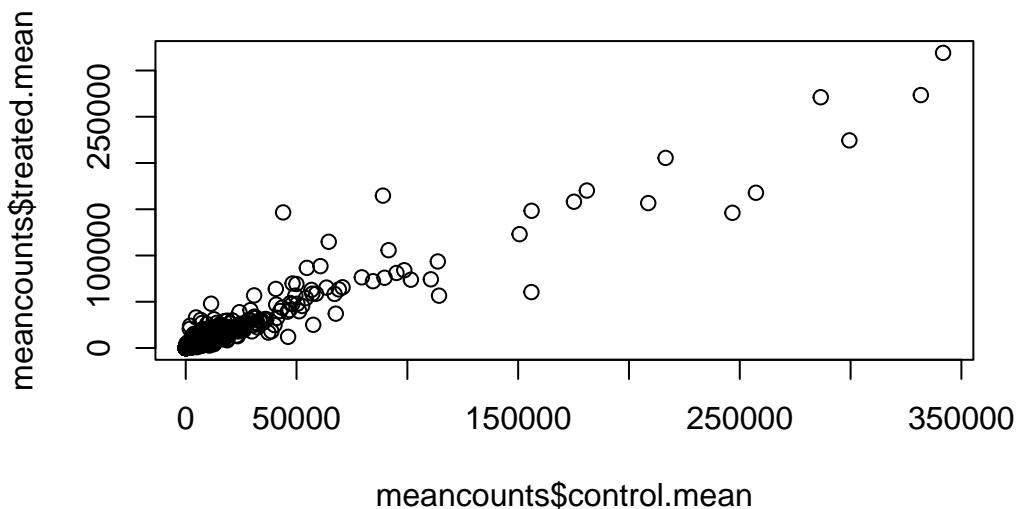
```
meancounts <- data.frame(control.mean, treated.mean)
head(meancounts)
```

	control.mean	treated.mean
ENSG000000000003	900.75	658.00
ENSG000000000005	0.00	0.00
ENSG000000000419	520.50	546.00

ENSG00000000457	339.75	316.50
ENSG00000000460	97.25	78.75
ENSG00000000938	0.75	0.00

Q5 (a). Create a scatter plot showing the mean of the treated samples against the mean of the control samples. Your plot should look something like the following.

```
plot(meancounts$control.mean, meancounts$treated.mean)
```



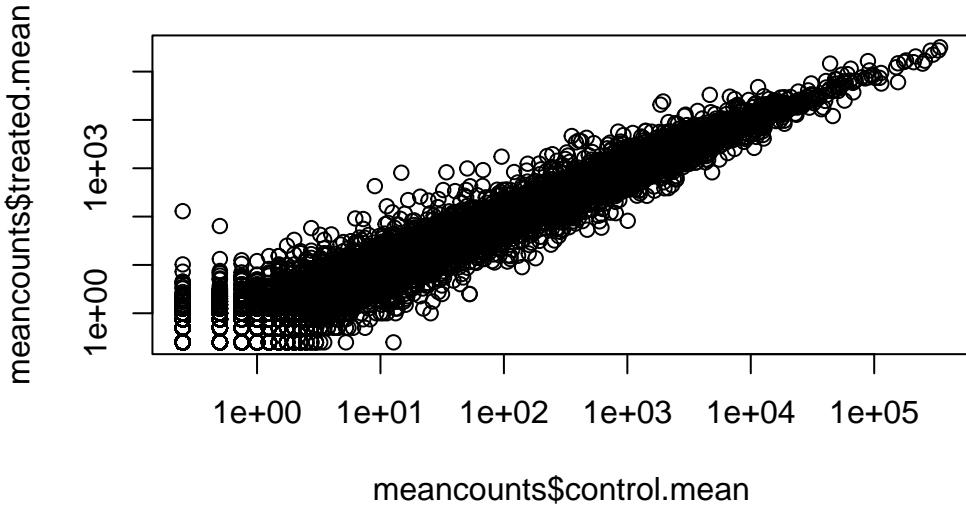
Q5 (b). You could also use the ggplot2 package to make this figure producing the plot below. What geom_?() function would you use for this plot? geom_point

Q6. Try plotting both axes on a log scale. What is the argument to plot() that allows you to do this? log()

```
plot(meancounts$control.mean, meancounts$treated.mean, log="xy")
```

Warning in xy.coords(x, y, xlabel, ylabel, log): 15032 x values <= 0 omitted from logarithmic plot

Warning in xy.coords(x, y, xlabel, ylabel, log): 15281 y values <= 0 omitted from logarithmic plot



Q7. What is the purpose of the arr.ind argument in the which() function call above? Why would we then take the first column of the output and need to call the unique() function? `arr.ind()` will return all the rows that do not have zeros in it. The `unique()` function will make sure that rows with zeros in both control and treated will not be counted twice.

```
zero.vals <- which(meancounts[, 1:2]==0, arr.ind=TRUE)

to.rm <- unique(zero.vals[, 1])
mycounts <- meancounts[-to.rm,]
head(mycounts)
```

	control.mean	treated.mean
ENSG00000000003	900.75	658.00
ENSG00000000419	520.50	546.00
ENSG00000000457	339.75	316.50
ENSG00000000460	97.25	78.75
ENSG00000000971	5219.00	6687.50
ENSG00000001036	2327.00	1785.75

The most useful and straightforward to understand is log2 transform.

```
#log2 of 1 =0 --> a value of 1 would mean that the treated and control were the same, so 1  
log2(20/20)
```

[1] 0

doubling

```
log2(40/20)
```

[1] 1

half the amount

```
log2(20/40)
```

[1] -1

4x the amount

```
log2(80/20)
```

[1] 2

add a log2 fold-change”

```
#fc means fold-change  
meancounts$log2fc <- log2(meancounts$treated.mean/meancounts$control.mean)
```

```
#NaN = "not a number"  
head(meancounts)
```

	control.mean	treated.mean	log2fc
ENSG000000000003	900.75	658.00	-0.45303916
ENSG000000000005	0.00	0.00	NaN
ENSG00000000419	520.50	546.00	0.06900279
ENSG00000000457	339.75	316.50	-0.10226805
ENSG00000000460	97.25	78.75	-0.30441833
ENSG00000000938	0.75	0.00	-Inf

Hmm.. we need to get rid of the genes where we have no count data as taking the log2 of these 0 counts does not tell us anything.

```
head(meancounts)
```

	control.mean	treated.mean	log2fc
ENSG000000000003	900.75	658.00	-0.45303916
ENSG000000000005	0.00	0.00	NaN
ENSG00000000419	520.50	546.00	0.06900279
ENSG00000000457	339.75	316.50	-0.10226805
ENSG00000000460	97.25	78.75	-0.30441833
ENSG00000000938	0.75	0.00	-Inf

```
head(meancounts$control.mean ==0)
```

```
[1] FALSE TRUE FALSE FALSE FALSE FALSE
```

```
head(meancounts$treated.mean ==0)
```

```
[1] FALSE TRUE FALSE FALSE FALSE TRUE
```

```
to.keep <- rowSums(meancounts[,1:2]==0) ==0  
mycounts <- meancounts[to.keep,]
```

Q8. Using the up.ind vector above can you determine how many up regulated genes we have at the greater than 2 fc level? 250 genes

```
sum(mycounts$log2fc > 2)
```

```
[1] 250
```

How many genes are upregulated at the log2fc level of +2 or greater?

```
sum(mycounts$log2fc >= 2)
```

```
[1] 314
```

and down regulated...

```
sum(mycounts$log2fc <= -2)
```

```
[1] 485
```

Q9. Using the down.ind vector above can you determine how many down regulated genes we have at the greater than 2 fc level? 367 genes

```
sum(mycounts$log2fc < (-2))
```

```
[1] 367
```

Q10. Do you trust these results? Why or why not? No because they exclude those equal to +2 and -2. They also do not really tell us how significant these results are compared to others in the sample.

We are missing the stats. are these big changes significant?

```
#DESeq2 analysis
```

```
library(DESeq2)
```

Like most bioconductor packages, DESeq wants its input and output in a specific format

```
#design is where in the metadata it tells you the design of the experiment
dds <- DESeqDataSetFromMatrix(countData=counts,
                                colData=metadata,
                                design=~dex)
```

```
converting counts to integer mode
```

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

```
dds
```

```

class: DESeqDataSet
dim: 38694 8
metadata(1): version
assays(1): counts
rownames(38694): ENSG00000000003 ENSG00000000005 ... ENSG00000283120
  ENSG00000283123
rowData names(0):
colnames(8): SRR1039508 SRR1039509 ... SRR1039520 SRR1039521
colData names(4): id dex celltype geo_id

```

The main DESeq function is called DESeq

```

  dds <- DESeq(dds)

estimating size factors

estimating dispersions

gene-wise dispersion estimates

mean-dispersion relationship

final dispersion estimates

fitting model and testing

  res <- results(dds)
  res

```

```

log2 fold change (MLE): dex treated vs control
Wald test p-value: dex treated vs control
DataFrame with 38694 rows and 6 columns
  baseMean log2FoldChange      lfcSE      stat     pvalue
  <numeric>      <numeric> <numeric> <numeric> <numeric>
ENSG00000000003   747.1942    -0.3507030  0.168246 -2.084470  0.0371175
ENSG00000000005    0.0000        NA        NA        NA        NA
ENSG000000000419   520.1342    0.2061078  0.101059  2.039475  0.0414026
ENSG000000000457   322.6648    0.0245269  0.145145  0.168982  0.8658106
ENSG000000000460   87.6826    -0.1471420  0.257007 -0.572521  0.5669691

```

```

...
...     ...
ENSG00000283115 0.000000      NA      NA      NA      NA
ENSG00000283116 0.000000      NA      NA      NA      NA
ENSG00000283119 0.000000      NA      NA      NA      NA
ENSG00000283120 0.974916 -0.668258 1.69456 -0.394354 0.693319
ENSG00000283123 0.000000      NA      NA      NA      NA
            padj
<numeric>
ENSG00000000003 0.163035
ENSG00000000005   NA
ENSG00000000419 0.176032
ENSG00000000457 0.961694
ENSG00000000460 0.815849
...
...
ENSG00000283115      NA
ENSG00000283116      NA
ENSG00000283119      NA
ENSG00000283120      NA
ENSG00000283123      NA

```

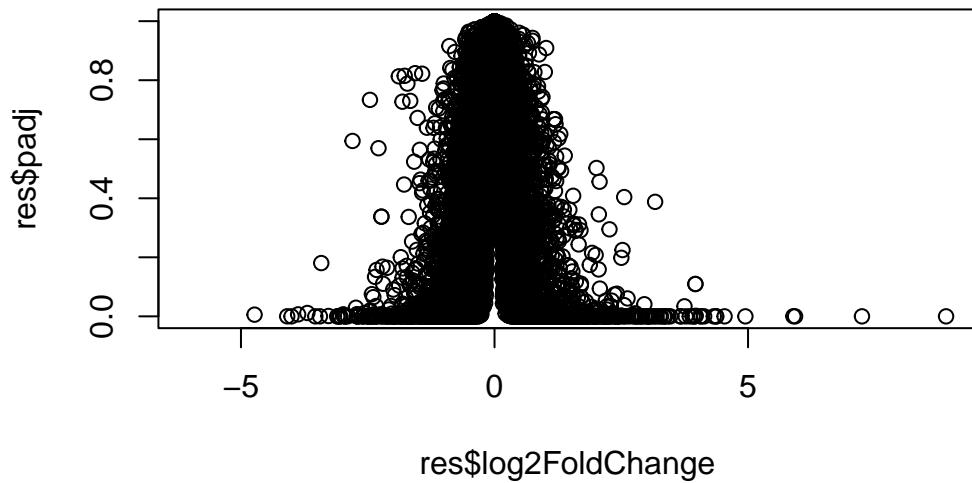
padj → adjusted p-value -there are so many different samples -scale by the amount of tests you do -can change things from significant to not significant compared to the rest of the data

Volcano Plots

A major summary figure of this type of analysis is called a volcano plot. The idea here is to keep our inner biologist and inner statistician happy with one cool plot!

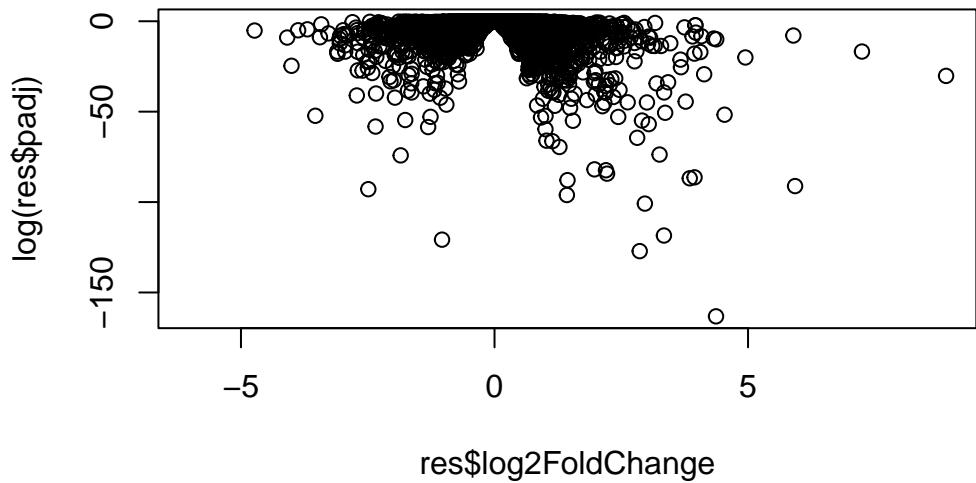
We want the more extreme values on the x axis (away from 0), and the lower p-values.

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



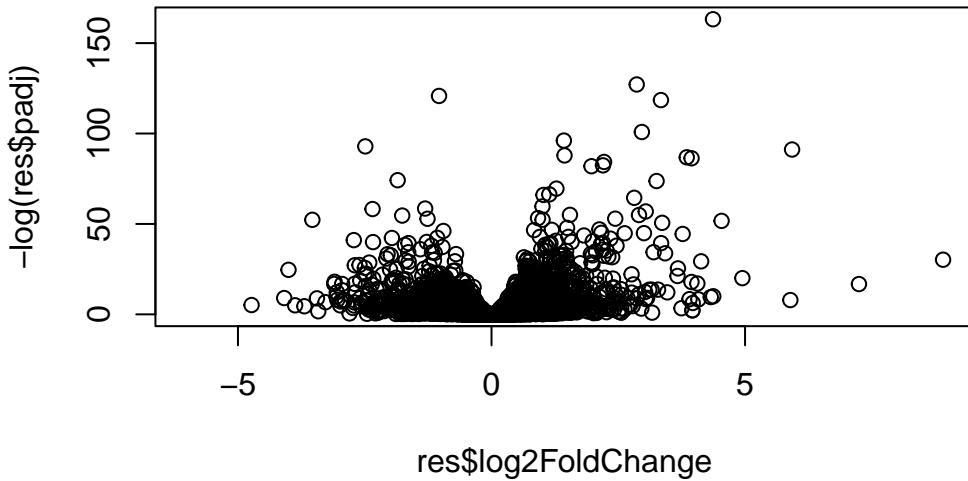
Improve this plot by taking the log of the p-value axis.

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



I want to flip this y-axis so the values I care about (i.e. the low p-values or high $\log(p\text{-value})$ are at the top).

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



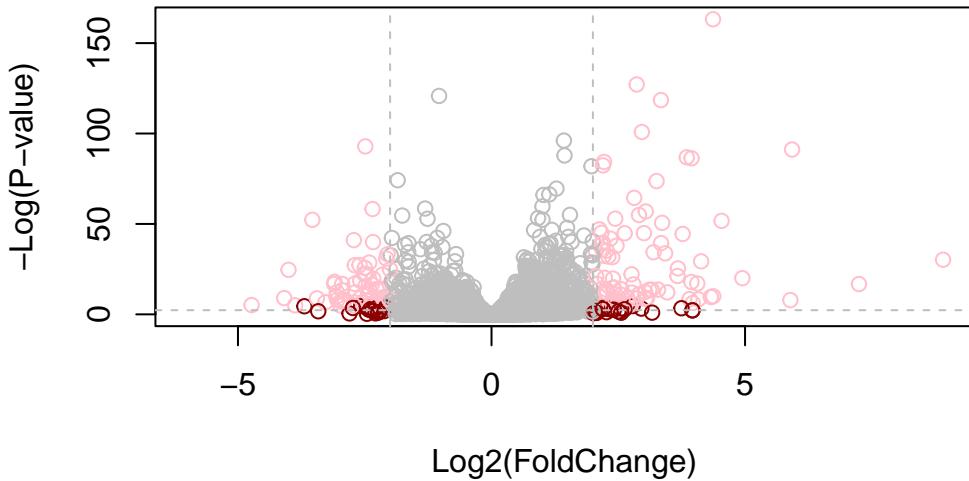
Let's finish up for today by adding some color to better highlight the subset of genes that we will focus on next day (i.e. those with big log2fc values at ± 2 threshold) and significant p-values (less than 0.05 for example).

```
# Setup our custom point color vector
mycols <- rep("gray", nrow(res))
mycols[ abs(res$log2FoldChange) > 2 ] <- "darkred"

inds <- (res$padj < 0.01) & (abs(res$log2FoldChange) > 2 )
mycols[ inds ] <- "pink"

# Volcano plot with custom colors
plot( res$log2FoldChange, -log(res$padj),
      col=mycols, ylab="-Log(P-value)", xlab="Log2(FoldChange)" )

# Cut-off lines
abline(v=c(-2,2), col="gray", lty=2)
abline(h=-log(0.1), col="gray", lty=2)
```



Adding annotation data

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

We will use one of Bioconductor's main annotation packages to help with mapping between various ID schemes. Here we load the `AnnotationDbi` package and the annotation data package for humans `org.Hs.eg.db`.

Look at what types of IFs I can translate between from the `org.Hs.eg.db` package with the `columns()` function.

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

```
[1] "ACCCNUM"          "ALIAS"           "ENSEMBL"          "ENSEMBLPROT"      "ENSEMLTRANS"
[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(x = org.Hs.eg.db,  
                      column= "SYMBOL",  
                      keys = rownames(res),  
                      keytype = "ENSEMBL")
```

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

Q11. Run the mapIds() function two more times to add the Entrez ID and UniProt accession and GENENAME as new columns called res\$entrez, res\$uniprot and res\$genename.

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

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

```
res$uniprot <- mapIds(org.Hs.eg.db,  
                      keys=row.names(res),  
                      column="UNIPROT",  
                      keytype="ENSEMBL",  
                      multiVals="first")
```

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

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

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

```
head(res)
```

```
log2 fold change (MLE): dex treated vs control
Wald test p-value: dex treated vs control
DataFrame with 6 rows and 10 columns
  baseMean log2FoldChange    lfcSE      stat   pvalue
  <numeric>     <numeric> <numeric> <numeric> <numeric>
ENSG000000000003 747.194195 -0.3507030 0.168246 -2.084470 0.0371175
ENSG000000000005 0.000000      NA        NA        NA        NA
ENSG00000000419 520.134160  0.2061078 0.101059  2.039475 0.0414026
ENSG00000000457 322.664844  0.0245269 0.145145  0.168982 0.8658106
ENSG00000000460 87.682625 -0.1471420 0.257007 -0.572521 0.5669691
ENSG00000000938 0.319167 -1.7322890 3.493601 -0.495846 0.6200029
  padj      symbol      entrez      uniprot
  <numeric> <character> <character> <character>
ENSG000000000003 0.163035    TSPAN6      7105 AOA024RCI0
ENSG000000000005  NA        TNMD       64102 Q9H2S6
ENSG00000000419 0.176032    DPM1       8813 060762
ENSG00000000457 0.961694    SCYL3      57147 Q8IZE3
ENSG00000000460 0.815849    C1orf112   55732 AOA024R922
ENSG00000000938  NA        FGR        2268 P09769
  genename
  <character>
ENSG000000000003 tetraspanin 6
ENSG000000000005 tenomodulin
ENSG00000000419 dolichyl-phosphate m..
ENSG00000000457 SCY1 like pseudokina..
ENSG00000000460 chromosome 1 open re..
ENSG00000000938 FGR proto-oncogene, ..
```

Pathway analysis

We will finish this lab with a quick pathway analysis. Here we play with just the **GAGE package**

```
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)
```

```
# Examine the first 2 pathways in this kegg set for humans  
head(kegg.sets.hs, 2)
```

```
$`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"
```

The main `gage()` function requires a named vector of fold changes where the names of the values are the Entrez gene IDs.

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

7105	64102	8813	57147	55732	2268	
-0.35070302		NA	0.20610777	0.02452695	-0.14714205	-1.73228897

Now,, let's run the pathway analysis

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

Now let's look at the object returned from gage()

```
attributes(keggres)
```

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

```
# Look at the first three down (less) pathways
head(keggres$less, 3)
```

		p.geomean	stat.mean	p.val
hsa05332	Graft-versus-host disease	0.0004250461	-3.473346	0.0004250461
hsa04940	Type I diabetes mellitus	0.0017820293	-3.002352	0.0017820293
hsa05310	Asthma	0.0020045888	-3.009050	0.0020045888
		q.val	set.size	exp1
hsa05332	Graft-versus-host disease	0.09053483	40	0.0004250461
hsa04940	Type I diabetes mellitus	0.14232581	42	0.0017820293
hsa05310	Asthma	0.14232581	29	0.0020045888

Let's pull up the highlighted pathways and show our differentially expressed genes on the pathway. I will use the 'hs' KEGG, I would get the pathway from KEGG and foldchange vector to show my genes

```
pathview(gene.data=foldchanges, pathway.id="hsa05310")
```

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

```
Info: Working in directory C:/Users/MacCh/OneDrive/Documents/BIMM 143/class12
```

```
Info: Writing image file hsa05310.pathview.png
```

Put this into my document

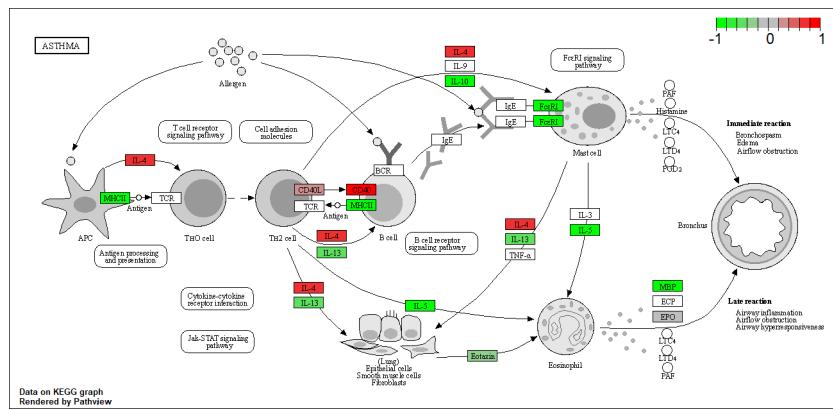


Figure 1: the Asthma pathway with my highlighted differentially expressed genes in order