

# lab12

Ariane

```
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
```

```
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, colAnys, 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, rowAnys, rowAvgsPerColSet,  
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
```

## Importing countData and colData

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

```
head(counts)
```

	SRR1039508	SRR1039509	SRR1039512	SRR1039513	SRR1039516
ENSG00000000003	723	486	904	445	1170
ENSG00000000005	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		
ENSG00000000003	1097	806	604		
ENSG00000000005	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?

There are 38694 genes in this dataset.

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

```
n.control <- sum(metadata$dex == "control")
```

There are 4 control cell lines in this dataset.

### **Extract and summarize the control samples**

To find out where the control samples are we need the metadata.

```

control <- metadata[metadata[, "dex"] == "control",]
control.counts <- counts[ , control$id]
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

```

control.mean <- rowMeans(control.counts)
head(control.mean)

```

ENSG000000000003	ENSG000000000005	ENSG000000000419	ENSG000000000457	ENSG000000000460
900.75	0.00	520.50	339.75	97.25
ENSG000000000938				
0.75				

```
library(dplyr)
```

Attaching package: 'dplyr'

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

combine

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

count

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

intersect, setdiff, union

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

intersect

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

collapse, desc, intersect, setdiff, slice, union

The following objects are masked from 'package:S4Vectors':

first, intersect, rename, setdiff, setequal, union

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

combine, intersect, setdiff, union

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

filter, lag

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

intersect, setdiff, setequal, union

```

control <- metadata %>% filter(dex=="control")
control.counts <- counts %>% select(control$id)
control.mean <- rowSums(control.counts)/4
head(control.mean)

```

ENSG000000000003	ENSG000000000005	ENSG000000000419	ENSG000000000457	ENSG00000000460
900.75	0.00	520.50	339.75	97.25
ENSG00000000938				
0.75				

Q3. How would you make the above code in either approach more robust?

We could make the above code in either approach more robust by not using `rowSums()` and manually dividing it by 4. Instead, we use `rowMeans()` which would give more robust results not matter the sample size.

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 <- metadata[metadata[, "dex"]=="treated",]
treated.counts <- counts[ ,treated$id]
treated.mean <- rowMeans(treated.counts)
head(treated.mean)

```

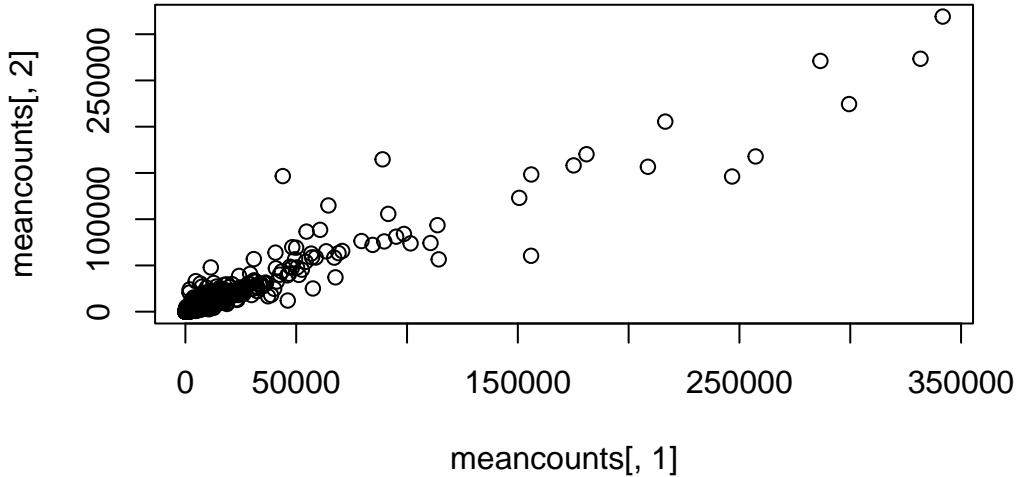
ENSG000000000003	ENSG000000000005	ENSG000000000419	ENSG000000000457	ENSG00000000460
658.00	0.00	546.00	316.50	78.75
ENSG00000000938				
0.00				

Store these results together in a new data frame called `meancounts`.

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

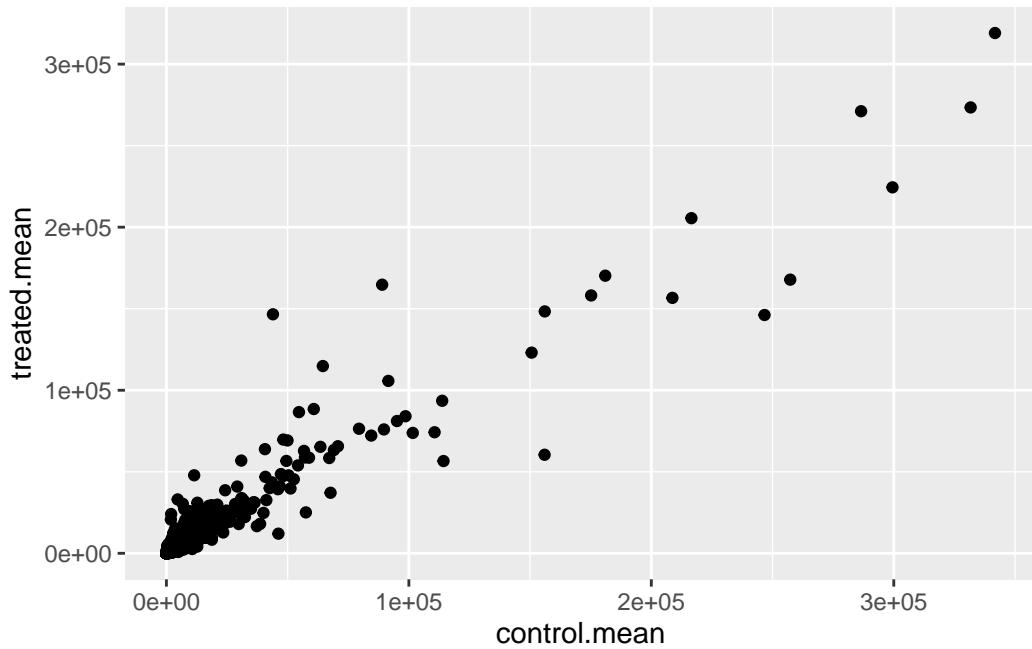
Q5 (a). Create a scatter plot showing the mean of the treated samples against the mean of the control samples. Your plot should look something like the following.

```
plot(meancounts[, 1], meancounts[, 2])
```



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?

```
library(ggplot2)
ggplot(meancounts, aes(control.mean, treated.mean)) + geom_point()
```

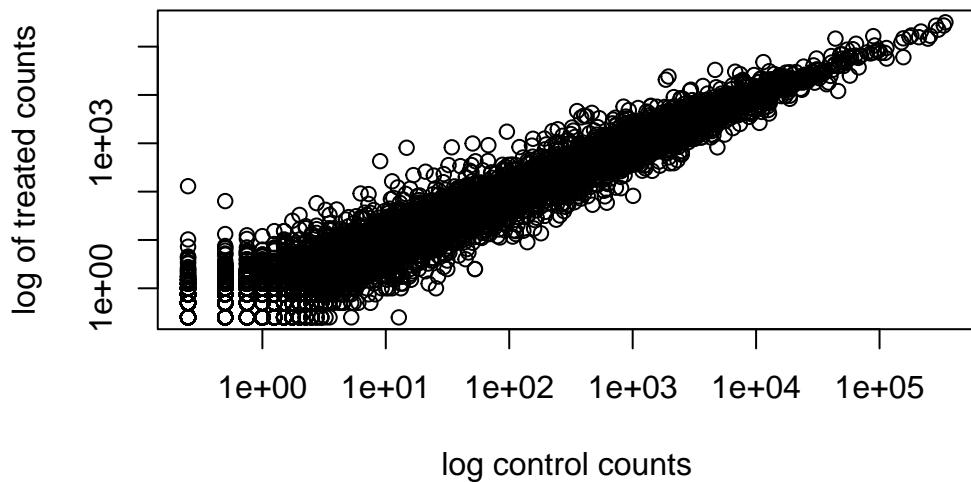


Q6. Try plotting both axes on a log scale. What is the argument to plot() that allows you to do this?

```
plot(meancounts[, 1], meancounts[,2], log="xy", xlab="log control counts", ylab ="log of t
```

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

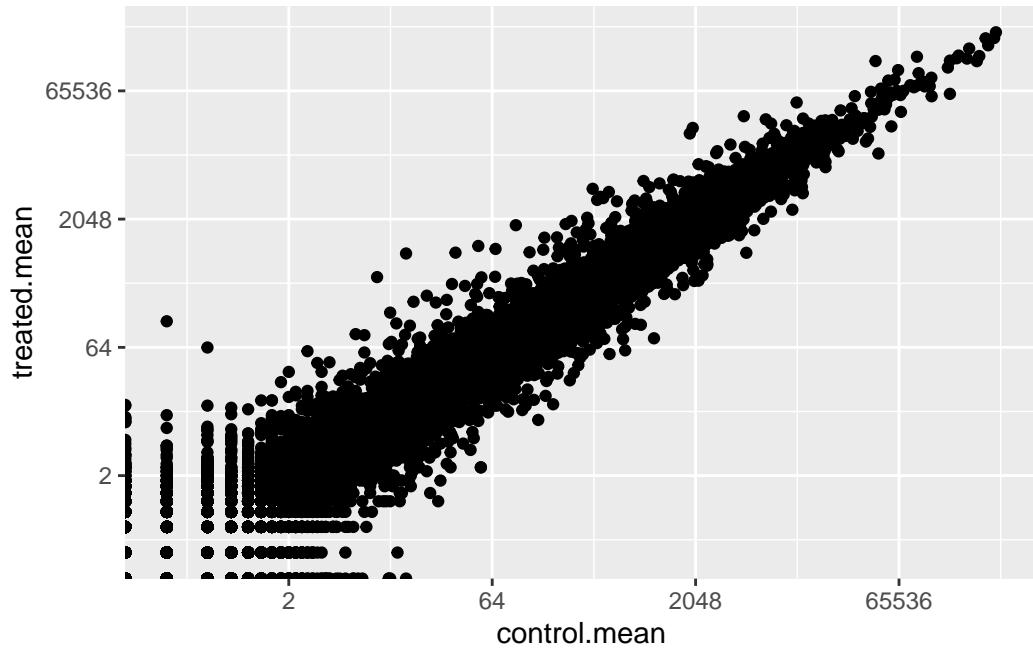
Warning in xy.coords(x, y, xlabel, ylabel, log): 15281 y values <= 0 omitted from logarithmic plot



```
ggplot(meancounts, aes(control.mean, treated.mean)) + geom_point() + scale_x_continuous(tr
```

Warning: Transformation introduced infinite values in continuous x-axis

Warning: Transformation introduced infinite values in continuous y-axis



This log2 transformation has the nice property where if there is no change the log2 value will be zero and if it doubles log2 value will be 1, and if halves it will be -1.

```
meancounts$log2fc <- log2(meancounts[, "treated.mean"] / meancounts[, "control.mean"])
head(meancounts)
```

	control.mean	treated.mean	log2fc
ENSG000000000003	900.75	658.00	-0.45303916
ENSG000000000005	0.00	0.00	NaN
ENSG000000000419	520.50	546.00	0.06900279
ENSG000000000457	339.75	316.50	-0.10226805
ENSG000000000460	97.25	78.75	-0.30441833
ENSG000000000938	0.75	0.00	-Inf

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?

```
head(which(meancounts[, 1:2]==0, arr.ind = TRUE))
```

	row	col
ENSG000000000005	2	1

```

ENSG00000004848 65 1
ENSG00000004948 70 1
ENSG00000005001 73 1
ENSG00000006059 121 1
ENSG00000006071 123 1

```

The arr.ind argument will give the names or index of the row where it exists a zero.

```

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	log2fc
ENSG00000000003	900.75	658.00	-0.45303916
ENSG00000000419	520.50	546.00	0.06900279
ENSG00000000457	339.75	316.50	-0.10226805
ENSG00000000460	97.25	78.75	-0.30441833
ENSG00000000971	5219.00	6687.50	0.35769358
ENSG00000001036	2327.00	1785.75	-0.38194109

We need to call `unique()` function because there might be multiple zeroes in one row.

A common threshold used for calling something differentially expressed is a log2(FoldChange) of greater than 2 or less than -2. Let's filter the dataset both ways to see how many genes are up or down-regulated.

```

up.ind <- mycounts$log2fc > 2
down.ind <- mycounts$log2fc < (-2)

```

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

```

sum(up.ind)

```

```

[1] 250

```

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

```
sum(down.ind)
```

```
[1] 367
```

Q10. Do you trust these results? Why or why not?

Not fully because we don't really know if these changes are statistically significant.

## DESeq2 Analysis

```
# load up DESeq2
library(DESeq2)
```

```
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 <- 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>
ENSG000000000003  747.1942   -0.3507030  0.168246 -2.084470 0.0371175
ENSG000000000005   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>
ENSG000000000003  0.163035
ENSG000000000005       NA
ENSG000000000419  0.176032
ENSG000000000457  0.961694
ENSG000000000460  0.815849
...
            ...
ENSG00000283115       NA
ENSG00000283116       NA
ENSG00000283119       NA
ENSG00000283120       NA
ENSG00000283123       NA

```

We can get some basic summary tallies using the `summary()` function.

```
summary(res, alpha=0.05)
```

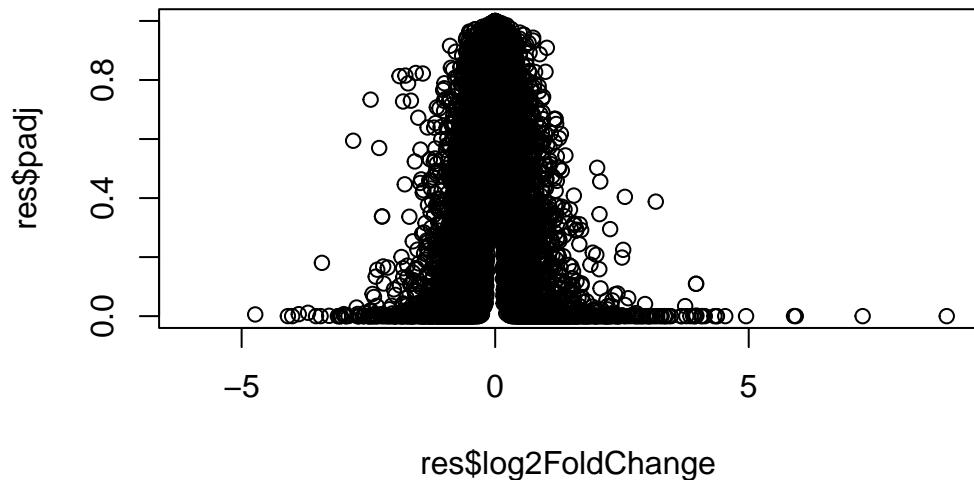
out of 25258 with nonzero total read count  
adjusted p-value < 0.05

```
LFC > 0 (up)      : 1242, 4.9%
LFC < 0 (down)    : 939, 3.7%
outliers [1]       : 142, 0.56%
low counts [2]     : 9971, 39%
(mean count < 10)
[1] see 'cooksCutoff' argument of ?results
[2] see 'independentFiltering' argument of ?results
```

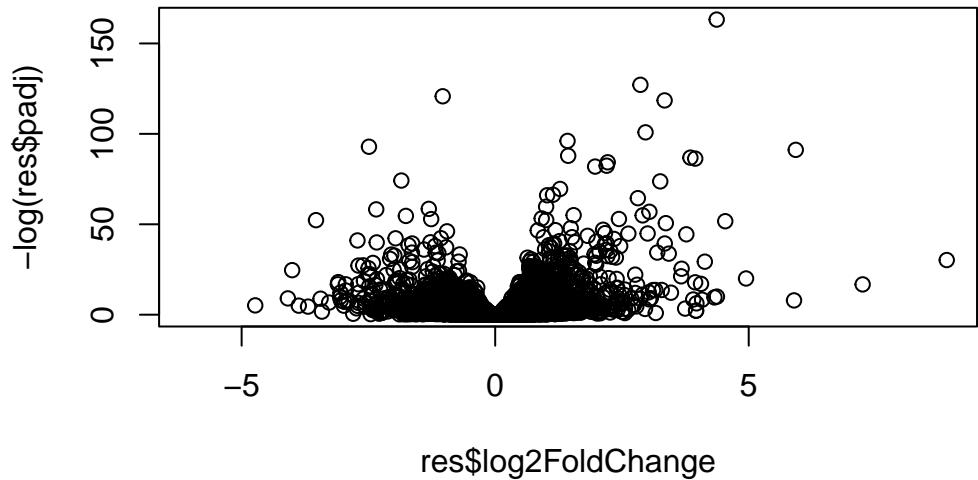
## Volcano plot

Make a summary plot of our results.

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



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



```

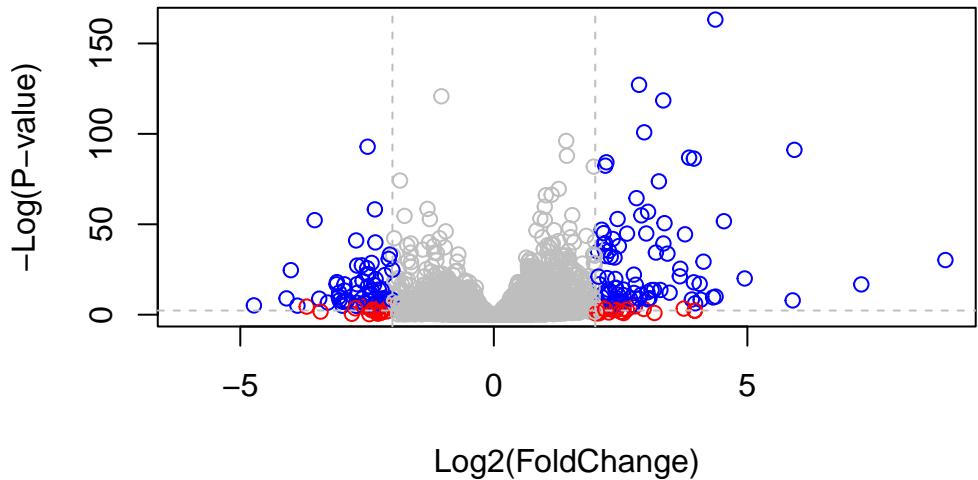
# Setup our custom point color vector
mycols <- rep("gray", nrow(res))
mycols[ abs(res$log2FoldChange) > 2 ] <- "red"

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

# 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)

```



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`.

```
library("AnnotationDbi")
```

```
Attaching package: 'AnnotationDbi'
```

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

```
select
```

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

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

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

```
[1] "ACCCNUM"      "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"
```

```
head(res)
```

```
log2 fold change (MLE): dex treated vs control
```

```
Wald test p-value: dex treated vs control
```

```
DataFrame with 6 rows and 6 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
ENSG000000000419	520.134160	0.2061078	0.101059	2.039475	0.0414026
ENSG000000000457	322.664844	0.0245269	0.145145	0.168982	0.8658106
ENSG000000000460	87.682625	-0.1471420	0.257007	-0.572521	0.5669691
ENSG000000000938	0.319167	-1.7322890	3.493601	-0.495846	0.6200029

	padj
	<numeric>
ENSG000000000003	0.163035
ENSG000000000005	NA
ENSG000000000419	0.176032
ENSG000000000457	0.961694
ENSG000000000460	0.815849
ENSG000000000938	NA

```
res$symbol <- mapIds(x=org.Hs.eg.db,
                      column="SYMBOL",
                      keys=rownames(res),
                      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 7 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
      <numeric> <character>
ENSG000000000003 0.163035    TSPAN6
ENSG000000000005  NA        TNMD
ENSG00000000419 0.176032    DPM1
ENSG00000000457 0.961694    SCYL3
ENSG00000000460 0.815849    C1orf112
ENSG00000000938  NA        FGR

res$entrez <- mapIds(x=org.Hs.eg.db,
                      column="ENTREZID",
                      keys=rownames(res),
                      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 8 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

```

```

ENSG000000000460 87.682625      -0.1471420 0.257007 -0.572521 0.5669691
ENSG000000000938 0.319167      -1.7322890 3.493601 -0.495846 0.6200029
                padj      symbol      entrez
                <numeric> <character> <character>
ENSG000000000003 0.163035      TSPAN6       7105
ENSG000000000005   NA          TNMD        64102
ENSG000000000419 0.176032      DPM1        8813
ENSG000000000457 0.961694      SCYL3       57147
ENSG000000000460 0.815849      C1orf112    55732
ENSG000000000938   NA          FGR         2268

```

## Pathway analysis

```
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"
```

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

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

```
attributes(keggres)
```

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

Now lets look at the object returned from `gage()`, i.e. our results here:

```
head(keggres$less,3)
```

		p.geomean	stat.mean
hsa04672	Intestinal immune network for IgA production	0.006043452	-2.560547
hsa04340	Hedgehog signaling pathway	0.013323955	-2.248547
hsa04916	Melanogenesis	0.030996739	-1.877209
		p.val	q.val
hsa04672	Intestinal immune network for IgA production	0.006043452	0.9763115
hsa04340	Hedgehog signaling pathway	0.013323955	0.9763115
hsa04916	Melanogenesis	0.030996739	0.9763115
		set.size	exp1
hsa04672	Intestinal immune network for IgA production	47	0.006043452
hsa04340	Hedgehog signaling pathway	56	0.013323955
hsa04916	Melanogenesis	101	0.030996739

```
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>).

```
#####
```

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

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

```
Info: Working in directory /Users/arianeyu/Desktop/BIMM 143/lab12
```

```
Info: Writing image file hsa04110.pathview.png
```

Put this into my document.

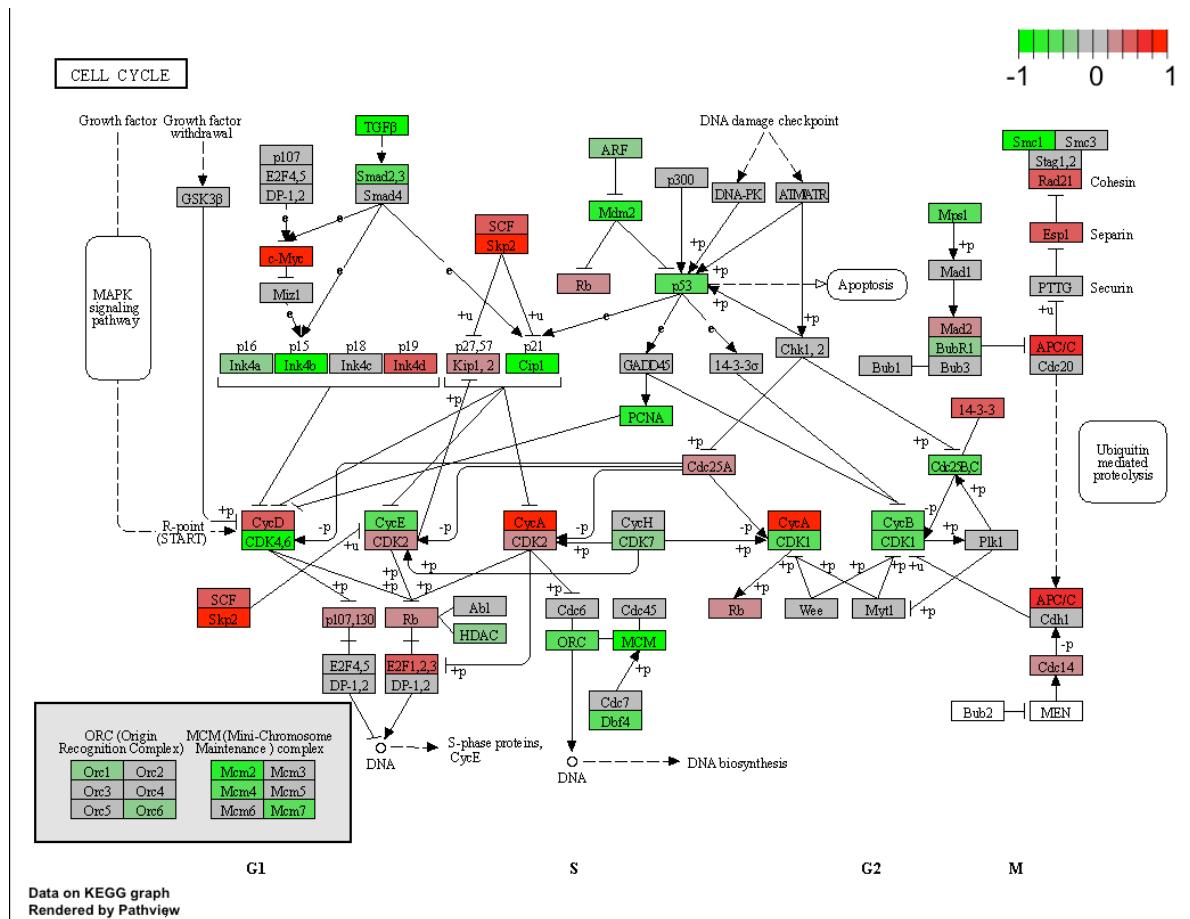


Figure 1: This Asthma pathway with my highlighted differentially expressed genes in color