

# Class 11: RNA-Seq continued

Hong (PID: A16558957)

2/22/2022

The data for this hands-on session comes from a published RNA-seq experiment where airway smooth muscle cells were treated with dexamethasone, a synthetic glucocorticoid steroid with anti-inflammatory effects

## Import/Read the data from Himes et al.

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

Have a little look see:

```
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
```

I always need to double check that the columns of my countdata and my coldata (metadata) match.

```
metadata$id
```

```
## [1] "SRR1039508" "SRR1039509" "SRR1039512" "SRR1039513" "SRR1039516"  
## [6] "SRR1039517" "SRR1039520" "SRR1039521"
```

```
colnames(counts)
```

```
## [1] "SRR1039508" "SRR1039509" "SRR1039512" "SRR1039513" "SRR1039516"  
## [6] "SRR1039517" "SRR1039520" "SRR1039521"
```

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

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

```
all(c(T,T,T))
```

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

I can use the `all()` function to make sure all my values match(e.e. all values are TRUE)

## 2. Extract control and treated counts for comparison

First lets extract the control counts columns

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

```
# Take the mean count value per gene (i.e. row)  
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
```

Now do the same thing for “treated” samples.

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

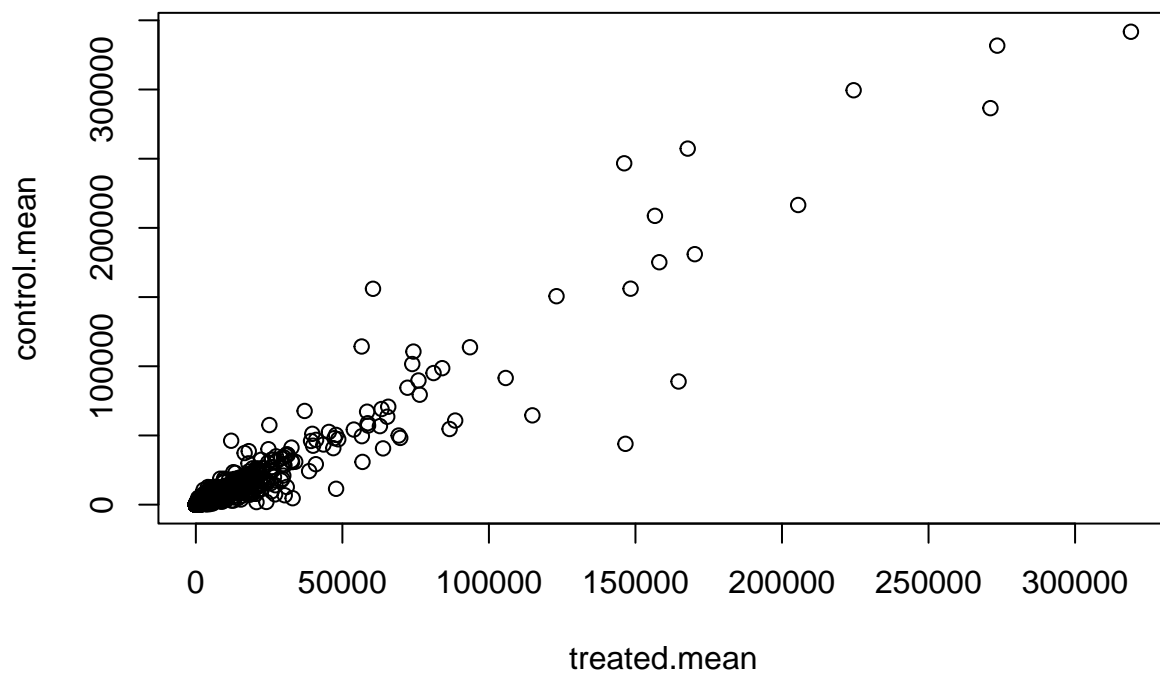
```
##           SRR1039509 SRR1039513 SRR1039517 SRR1039521
## ENSG00000000003      486      445      1097      604
## ENSG00000000005        0        0        0        0
## ENSG00000000419      523      371      781      509
## ENSG00000000457      258      237      447      324
## ENSG00000000460       81       66       94       74
## ENSG00000000938        0        0        0        0
```

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

```
## ENSG00000000003 ENSG00000000005 ENSG00000000419 ENSG00000000457 ENSG00000000460
##           658.00           0.00          546.00          316.50          78.75
## ENSG00000000938
##           0.00
```

Now we can make a plot comparing treated vs control

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

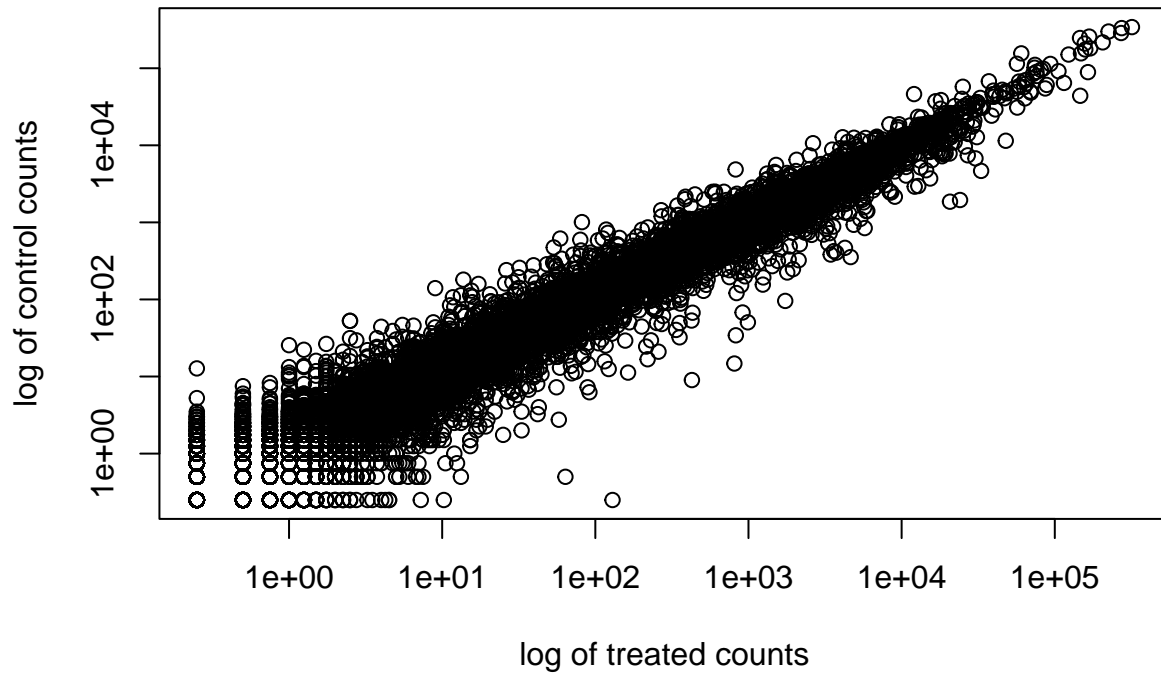


When we see data that is so skewed like this over quite a wide range of values we start to think of log transformations to make our analysis easier.

```
plot(treated.mean, control.mean, log = "xy",  
     xlab = "log of treated counts",  
     ylab = "log of control counts")
```

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

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



We are after changes in gene expression: treated vs control and this would represent points (i.e. genes) that do not lie on the diagonal.

We like to work with log 2 values

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

```
## [1] 1
```

```
log2(640/20)
```

```
## [1] 5
```

Now let's calculate the log2 fold change

```
log2fc <- log2(treated.mean/control.mean)
```

Store my work so far

```
meancounts <- data.frame(control.mean, treated.mean, log2fc)
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
```

Filter our data to remove genes with zero expression values.

```
c(10,200,0,40) == 0
```

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

```
which(c(10,200,0,40) == 0)
```

```
## [1] 3
```

```
z <- data.frame(x=c(10,0,30,40),
                y=c(10,0,30,0))
which(z == 0, arr.ind = TRUE)
```

```
##      row col
## [1,]   2   1
## [2,]   2   2
## [3,]   4   2
```

```
i <- which(z == 0, arr.ind = TRUE)
unique(i[,1])
```

```
## [1] 2 4
```

Now do it for our real data set

```
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
## ENSG000000000003      900.75      658.00 -0.45303916
## ENSG000000000419      520.50      546.00  0.06900279
## ENSG000000000457      339.75      316.50 -0.10226805
## ENSG000000000460       97.25       78.75 -0.30441833
## ENSG000000000971     5219.00     6687.50  0.35769358
## ENSG00000001036     2327.00     1785.75 -0.38194109
```

How many genes do we have left?

```
nrow(mycounts)
```

```
## [1] 21817
```

A common threshold used for calling something differentially expressed is a  $\log_2(\text{FoldChange})$  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" genes ...

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

```
## [1] 250
```

"Down" genes ...

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

```
## [1] 367
```

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

## Loading required package: Biobase

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

##
## Attaching package: 'Biobase'

```

```
## The following object is masked from 'package:MatrixGenerics':  
##  
## rowMedians
```

```
## The following objects are masked from 'package:matrixStats':  
##  
## anyMissing, rowMedians
```

```
citation("DESeq2")
```

```
##  
## Love, M.I., Huber, W., Anders, S. Moderated estimation of fold change  
## and dispersion for RNA-seq data with DESeq2 Genome Biology 15(12):550  
## (2014)  
##  
## A BibTeX entry for LaTeX users is  
##  
## @Article{,  
## title = {Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2},  
## author = {Michael I. Love and Wolfgang Huber and Simon Anders},  
## year = {2014},  
## journal = {Genome Biology},  
## doi = {10.1186/s13059-014-0550-8},  
## volume = {15},  
## issue = {12},  
## pages = {550},  
## }
```

```
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>
## 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				
## ENSG000000000419	0.176032				
## ENSG000000000457	0.961694				
## ENSG000000000460	0.815849				
## ...	...				
## ENSG00000283115	NA				
## ENSG00000283116	NA				
## ENSG00000283119	NA				
## ENSG00000283120	NA				
## ENSG00000283123	NA				

```
summary(res)
```

```
##
## out of 25258 with nonzero total read count
## adjusted p-value < 0.1
## LFC > 0 (up)      : 1563, 6.2%
## LFC < 0 (down)    : 1188, 4.7%
## 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
```

```
res05 <- results(dds, alpha=0.05)
summary(res05)
```

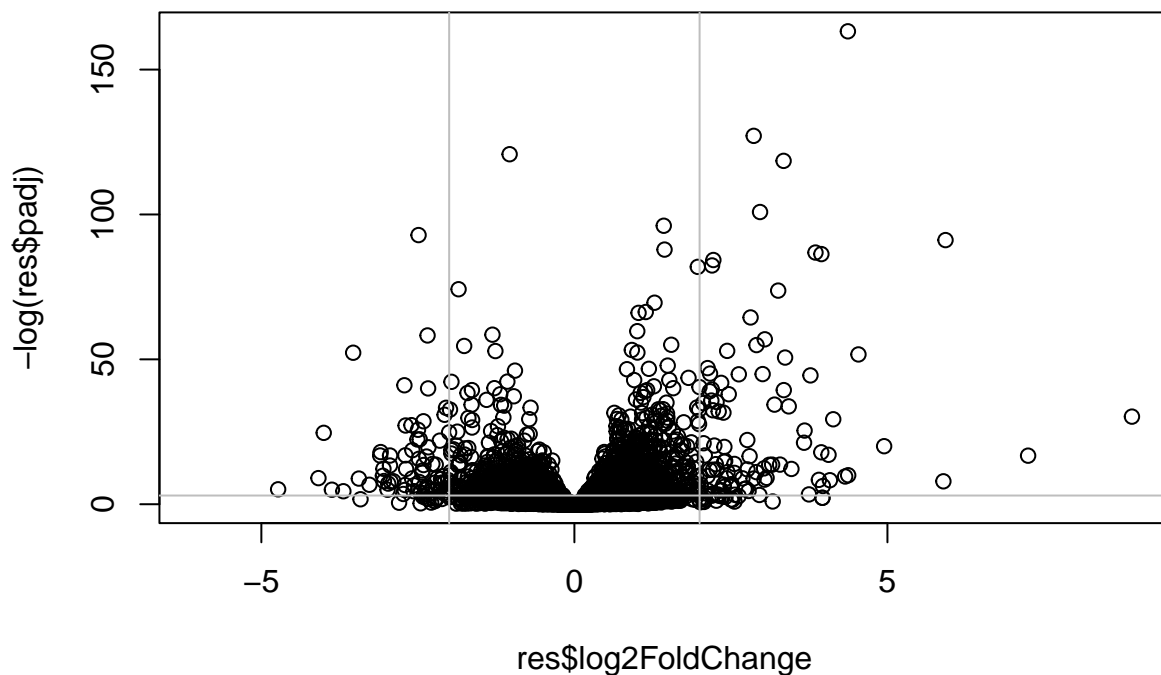
```
##
## out of 25258 with nonzero total read count
## adjusted p-value < 0.05
```

```
## LFC > 0 (up)      : 1236, 4.9%
## LFC < 0 (down)    : 933, 3.7%
## outliers [1]      : 142, 0.56%
## low counts [2]     : 9033, 36%
## (mean count < 6)
## [1] see 'cooksCutoff' argument of ?results
## [2] see 'independentFiltering' argument of ?results
```

## Volcano plots

Let's make a commonly produced visualization from this data, namely a so-called Volcano plot. These summary figures are frequently used to highlight the proportion of genes that are both significantly regulated and display a high fold change.

```
plot(res$log2FoldChange, -log(res$padj))
abline(h=-log(0.05), col = "gray")
abline(v=c(-2,2), col="gray")
```



I want to polish this main results figure by adding color to the genes i will focus on next day

```
# I will start by making a gray vector for everything
mycols <- rep("gray", nrow(res))

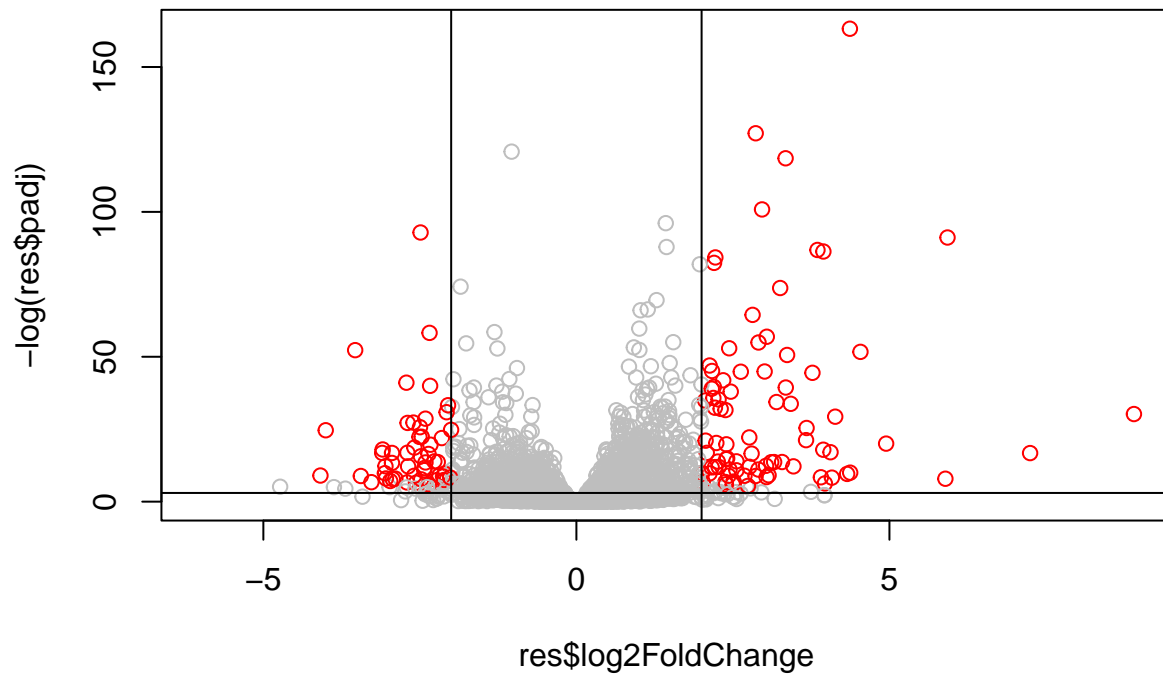
# Now I will overwrite the small padj values
mycols[res$padj < 0.005] <- "red"
```

```

# Now if my log2foldchange is small I will make them gray
mycols[abs(res$log2FoldChange) < 2] <- "gray"

plot(res$log2FoldChange, -log(res$padj), col = mycols)
abline(h=-log(0.05))
abline(v=c(-2,2))

```



This is a common overall summary figure because it combines big changes (in terms of log2 foldchange) and significant changes (in terms of p-value) all in one figure.