# R          DESeq

RNA                                                                      DESeq

## 1.

### DESeq

> source("http://bioconductor.org/biocLite.R")

> biocLite("DESeq")

```
package 'DBI' successfully unpacked and MD5 sums checked
package 'RSQLite' successfully unpacked and MD5 sums checked
package 'IRanges' successfully unpacked and MD5 sums checked
package 'xtable' successfully unpacked and MD5 sums checked
package 'XML' successfully unpacked and MD5 sums checked
package 'BiocGenerics' successfully unpacked and MD5 sums checked
package 'AnnotationDbi' successfully unpacked and MD5 sums checked
package 'annotate' successfully unpacked and MD5 sums checked
package 'Biobase' successfully unpacked and MD5 sums checked
package 'locfit' successfully unpacked and MD5 sums checked
package 'genefilter' successfully unpacked and MD5 sums checked
package 'geneplotter' successfully unpacked and MD5 sums checked
package 'RColorBrewer' successfully unpacked and MD5 sums checked
package 'DESeq' successfully unpacked and MD5 sums checked
```

```
package 'cluster' successfully unpacked and MD5 sums checked
package 'deSolve' successfully unpacked and MD5 sums checked
package 'foreign' successfully unpacked and MD5 sums checked
package 'KernSmooth' successfully unpacked and MD5 sums checked
package 'lattice' successfully unpacked and MD5 sums checked
package 'Matrix' successfully unpacked and MD5 sums checked
package 'mgcv' successfully unpacked and MD5 sums checked
package 'nnet' successfully unpacked and MD5 sums checked
package 'plotrix' successfully unpacked and MD5 sums checked
package 'rpart' successfully unpacked and MD5 sums checked
package 'survival' successfully unpacked and MD5 sums checked
```

pasilla.

## 2.

### 2.1

i          j                    j                    i                         reads

pasilla 　　　　　　 system.file

```
> datafile = system.file( "extdata/pasilla_gene_counts.tsv", package="pasilla" )
> datafile
[1] "D:/Program Files/R/R-2.15.3/library/pasilla/extdata/pasilla_gene_counts.tsv"
```

R 　　　　　　　　　 read.table

```
> pasillaCountTable = read.table( datafile, header=TRUE, row.names=1 )
> head( pasillaCountTable )
```

| | untreated1 | untreated2 | untreated3 | untreated4 | treated1 | treated2 | treated3 |
|---|---|---|---|---|---|---|---|
| FBgn0000003 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| FBgn0000008 | 92 | 161 | 76 | 70 | 140 | 88 | 70 |
| FBgn0000014 | 5 | 1 | 0 | 0 | 4 | 0 | 0 |
| FBgn0000015 | 0 | 2 | 1 | 2 | 1 | 0 | 0 |
| FBgn0000017 | 4664 | 8714 | 3564 | 3150 | 6205 | 3072 | 3334 |
| FBgn0000018 | 583 | 761 | 245 | 310 | 722 | 299 | 308 |

## 2.2

( ) 　　 ( )

data.frame 　　　　　　 7

```
> pasillaDesign = data.frame(
+ row.names = colnames( pasillaCountTable ),
+ condition = c( "untreated", "untreated", "untreated",
+ "untreated", "treated", "treated", "treated" ),
+ libType = c( "single-end", "single-end", "paired-end",
+ "paired-end", "single-end", "paired-end", "paired-end" ) )
> pasillaDesign
```

| | condition | libType |
|---|---|---|
| untreated1 | untreated | single-end |
| untreated2 | untreated | single-end |
| untreated3 | untreated | paired-end |
| untreated4 | untreated | paired-end |
| treated1 | treated | single-end |
| treated2 | treated | paired-end |
| treated3 | treated | paired-end |

R

single-end 　　 paired-end 　　　　　　　　　 paired-end

```
> pairedSamples = pasillaDesign$libType == "paired-end"
> countTable = pasillaCountTable[ , pairedSamples ]
> condition = pasillaDesign$condition[ pairedSamples ]

> head(countTable)
```

```
          untreated3 untreated4 treated2 treated3
FBgn0000003         0          0        0        1
FBgn0000008        76         70       88       70
FBgn0000014         0          0        0        0
FBgn0000015         1          2        0        0
FBgn0000017      3564       3150     3072     3334
FBgn0000018       245        310      299      308
```

> condition

```
[1] untreated untreated treated   treated
Levels: treated untreated
```

> #not run

> condition = factor( c( "untreated", "untreated", "treated", "treated" ) )

     CountDataSet  DESeq

> library( "DESeq" )

> cds = newCountDataSet( countTable, condition )

## 2.3

   estimateSizeFactors

> cds = estimateSizeFactors( cds )

> sizeFactors( cds )

```
untreated3 untreated4   treated2   treated3
 0.8730966  1.0106112  1.0224517  1.1145888
```

   counts

> head( counts( cds, normalized=TRUE ) )

```
          untreated3 untreated4   treated2   treated3
FBgn0000003   0.000000    0.00000    0.00000   0.8971919
FBgn0000008  87.046493   69.26502   86.06763  62.8034302
FBgn0000014   0.000000    0.00000    0.00000   0.0000000
FBgn0000015   1.145349    1.97900    0.00000   0.0000000
FBgn0000017 4082.022370 3116.92579 3004.54278 2991.2376629
FBgn0000018  280.610404  306.74508  292.43434  276.3350930
```

## 3.

DESeq

> cds = estimateDispersions( cds )

   estimateDispersions

    fitInfo

> str( fitInfo(cds) )

```
List of 5
 $ perGeneDispEsts: num [1:14599] -0.4696 0.0237 NaN -0.9987 0.0211 ...
 $ dispFunc       :function (q)
  ..- attr(*, "coefficients")= Named num [1:2] 0.00524 1.16816
  .. ..- attr(*, "names")= chr [1:2] "asymptDisp" "extraPois"
  ..- attr(*, "fitType")= chr "parametric"
 $ fittedDispEsts : num [1:14599] 5.21332 0.02055 Inf 1.5008 0.00559 ...
 $ df             : int 2
 $ sharingMode    : chr "maximum"
```
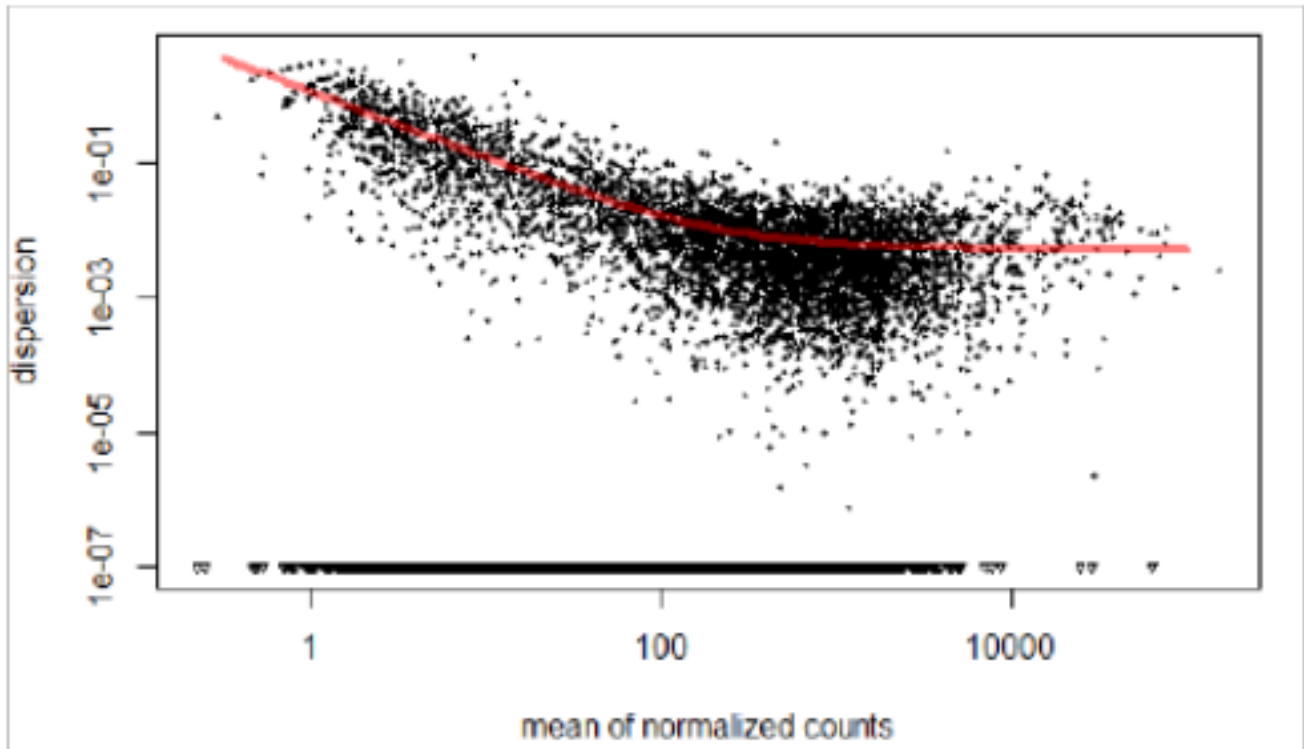
plotDispEsts

`> plotDispEsts( cds )`



1        (    )        (    )

cds

`> head( fData(cds) )`

```
           disp_pooled
FBgn0000003 5.213318025
FBgn0000008 0.023683636
FBgn0000014         Inf
FBgn0000015 1.500797095
FBgn0000017 0.021099504
FBgn0000018 0.009280402
```

## 4.

### 4.1

"untreated" "treated"                                nbinomTest

`> res = nbinomTest( cds, "untreated", "treated" )`

`> head(res)`

| | id | baseMean | baseMeanA | baseMeanB | foldChange | log2FoldChange | pval | padj |
|---|---|---|---|---|---|---|---|---|
| 1 | FBgn0000003 | 0.2242980 | 0.000000 | 0.4485959 | Inf | Inf | 1.0000000 | 1.0000000 |
| 2 | FBgn0000008 | 76.2956431 | 78.155755 | 74.4355310 | 0.9523999 | -0.07036067 | 0.8354725 | 1.0000000 |
| 3 | FBgn0000014 | 0.0000000 | 0.000000 | 0.0000000 | NaN | NaN | NA | NA |
| 4 | FBgn0000015 | 0.7810873 | 1.562175 | 0.0000000 | 0.0000000 | -Inf | 0.4160556 | 1.0000000 |
| 5 | FBgn0000017 | 3298.6821506 | 3599.474076 | 2997.8902236 | 0.8328690 | -0.26383857 | 0.2414208 | 0.8811746 |
| 6 | FBgn0000018 | 289.0312286 | 293.677741 | 284.3847165 | 0.9683564 | -0.04638999 | 0.7572819 | 1.0000000 |

| | |
|---|---|
| id | feature identier |
| baseMean | mean normalised counts, averaged over all samples from both conditions |
| baseMeanA | mean normalised counts from condition A |
| baseMeanB | mean normalised counts from condition B |
| foldChange | fold change from condition A to B |
| log2FoldChange | the logarithm (to basis 2) of the fold change |
| pval | pvalue for the statistical signicance of this change |
| padj | pvalue adjusted for multiple testing with the Benjamini-Hochberg procedure (see the R functionp.adjust), which controls false discovery rate (FDR) |

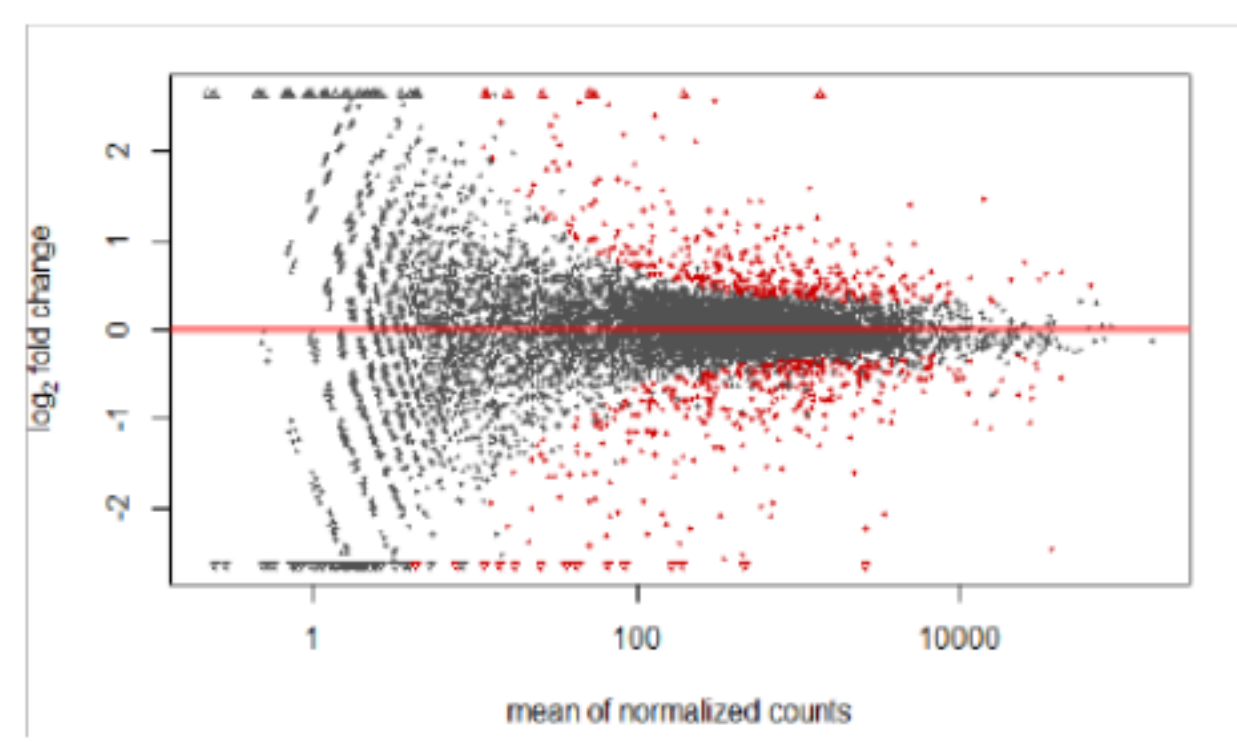log2                                                         10%FDR

> plotMA(res)



2 log2
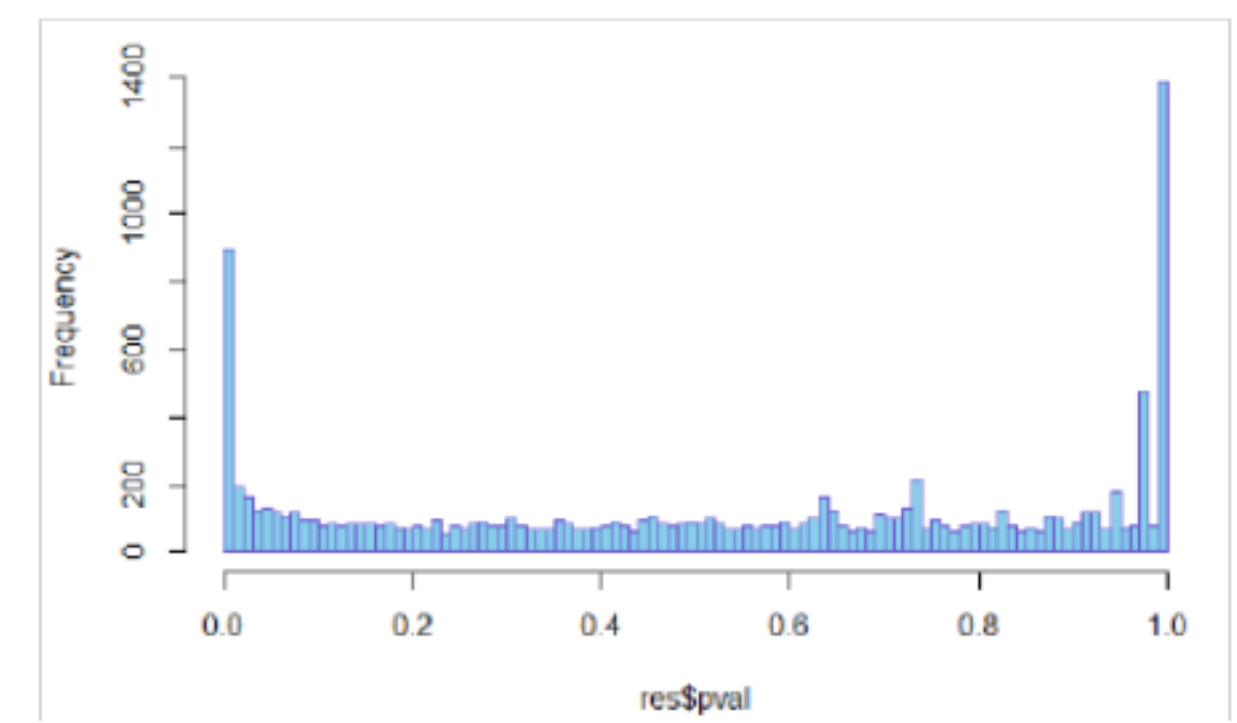


3      nbinomTest      p

> hist(res$pval, breaks=100, col="skyblue", border="slateblue", main="")

FDR

> resSig = res[ res$padj < 0.1, ]


> head( resSig[ order(resSig$pval), ] )

```
         id  baseMean baseMeanA baseMeanB foldChange log2FoldChange         pval         padj
9831  FBgn0039155  463.4369  884.9640   41.90977  0.0473576      -4.400260 1.641210e-124 1.887556e-120
2366  FBgn0025111 1340.2282  311.1697 2369.28680  7.6141316       2.928680 3.496915e-107 2.010901e-103
612   FBgn0003360 2544.2512 4513.9457  574.55683  0.1272848      -2.973868 1.552884e-99  5.953239e-96
5192  FBgn0029167 2551.3113 4210.9571  891.66551  0.2117489      -2.239574 4.346335e-78  1.249680e-74
10305 FBgn0039827  188.5927  357.3299   19.85557  0.0555665      -4.169641 1.189136e-65  2.735251e-62
6948  FBgn0035085  447.2485  761.1898  133.30718  0.1751300      -2.513502 3.145997e-56  6.030352e-53
```

> head( resSig[ order( resSig$foldChange, -resSig$baseMean ), ] )

```
         id   baseMean baseMeanA baseMeanB foldChange log2FoldChange         pval         padj
14078 FBgn0259236   4.269698  8.539395  0.000000 0.00000000          -Inf 1.305127e-03  2.601433e-02
13584 FBgn0085359  36.471017 71.026377  1.915658 0.02697108     -5.212443 2.404195e-09  2.194496e-07
9831  FBgn0039155 463.436895 884.964018 41.909773 0.04735760     -4.400260 1.641210e-124 1.887556e-120
2277  FBgn0024288  42.564985 80.890483  4.239487 0.05241021     -4.254008 5.568809e-20  1.940814e-17
10305 FBgn0039827 188.592732 357.329894 19.855571 0.05556650     -4.169641 1.189136e-65  2.735251e-62
6495  FBgn0034434  82.886729 155.091929 10.681529 0.06887224     -3.859934 5.936541e-32  5.252012e-29
```

> head( resSig[ order( -resSig$foldChange, -resSig$baseMean ), ] )

```
         id   baseMean baseMeanA  baseMeanB foldChange log2FoldChange         pval         padj
6030  FBgn0033764  53.94641 10.257418   97.63340  9.518321       3.250707 9.947141e-18  2.933386e-15
13079 FBgn0063667  11.77349  2.551675   20.99531  8.228051       3.040551 1.433352e-04  4.222521e-03
7020  FBgn0035189 197.46886 45.299371  349.63836  7.718393       2.948301 6.433138e-15  1.345228e-12
5499  FBgn0037290  50.45147 11.663744   89.23920  7.650990       2.935647 1.035267e-14  2.052865e-12
2366  FBgn0025111 1340.22823 311.169666 2369.28680  7.614132       2.928680 3.496915e-107 2.010901e-103
7264  FBgn0035539  15.99844  4.191774   27.80510  6.633254       2.729717 6.759276e-05  2.141555e-03
```

                    CSV

> write.csv( res, file="My Pasilla Analysis Result Table.csv" )

                                        untreated

> ncu = counts( cds, normalized=TRUE )[ , conditions(cds)=="untreated" ]

Ncu

>plotMA(data.frame(baseMean = rowMeans(ncu),

+       log2FoldChange = log2( ncu[,2] / ncu[,1] )),

+       col = "black")

        4



        4       untreated              log2


4.2


> cdsUUT = cds[ , 1:3]

```
> pData( cdsUUT )
```

```
         sizeFactor condition
untreated3  0.8730966 untreated
untreated4  1.0106112 untreated
treated2    1.0224517   treated
```
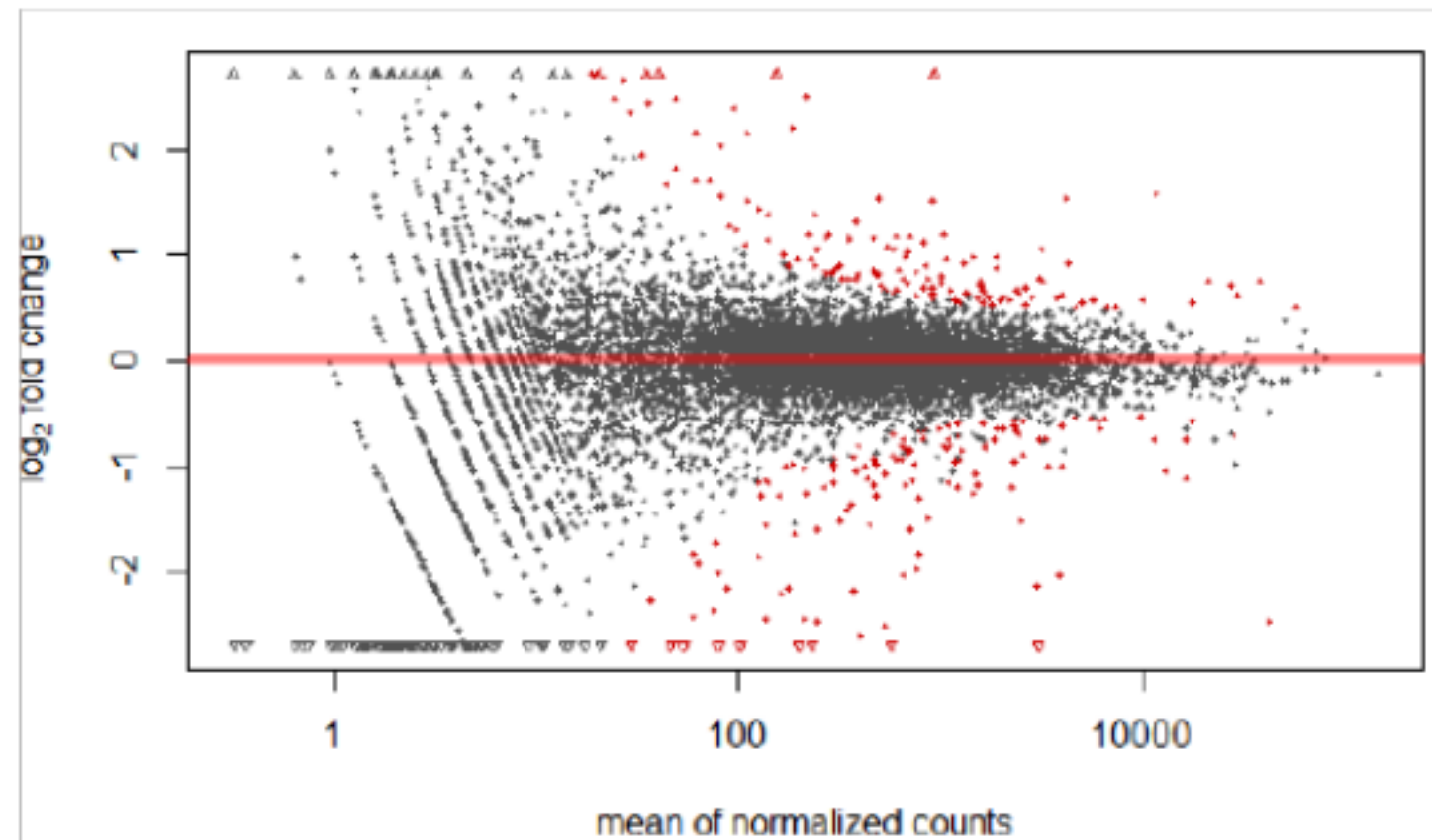


5 " treated" vs. untreated"        MvA    ,        treated        untreated

```
>cdsUUT = estimateSizeFactors( cdsUUT )
>cdsUUT = estimateDispersions( cdsUUT )
>resUUT = nbinomTest( cdsUUT, "untreated", "treated" )

>plotMA(resUUT)
```

## 4.3

DESeq
```
> cds2 = cds[ ,c( "untreated3", "treated3" ) ]
> cds2 = estimateDispersions( cds2, method="blind", sharingMode="fit-only" )
> res2 = nbinomTest( cds2, "untreated", "treated" )
> plotMA(res2)
```
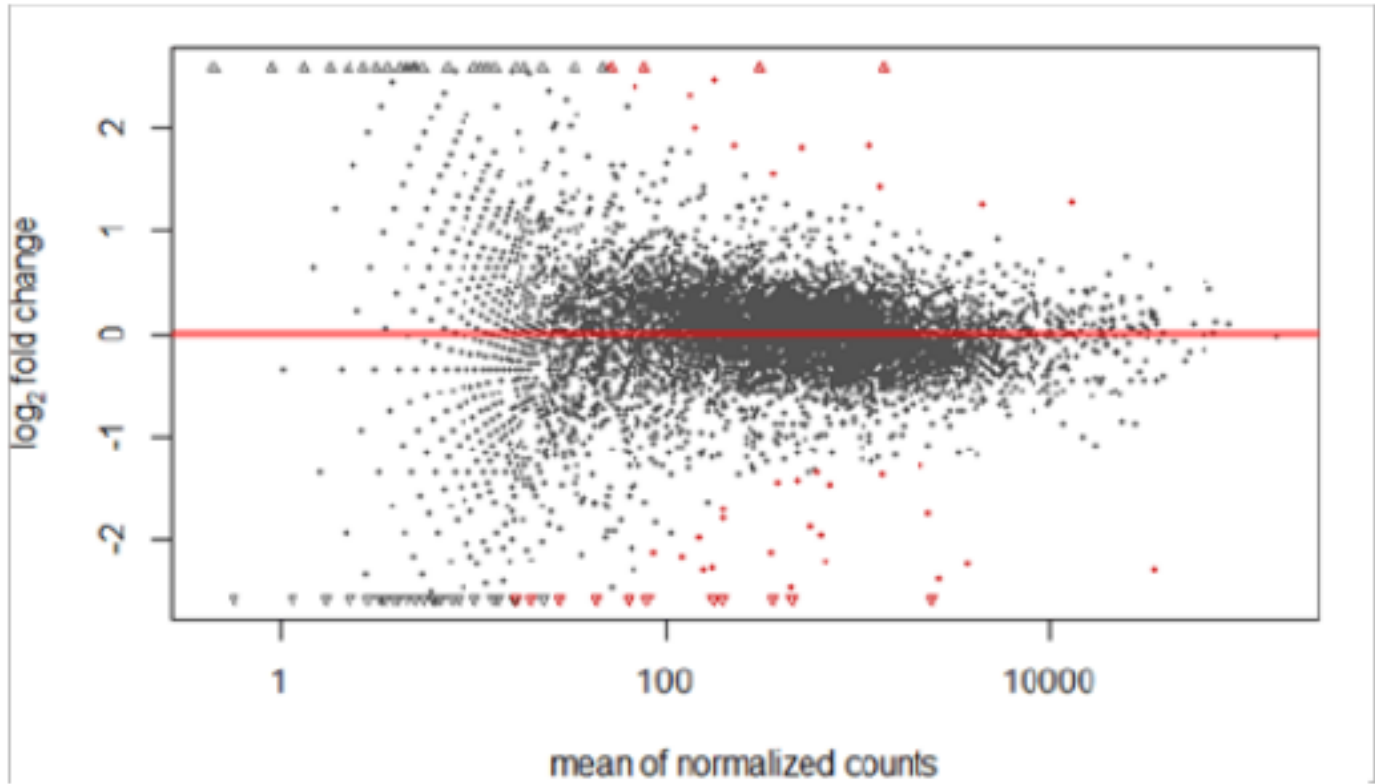
6 MvA

> addmargins( table( res_sig = res$padj < .1, res2_sig = res2$padj < .1 ) )

```
        res2_sig
res_sig FALSE  TRUE   Sum
  FALSE 10084     1 10085
   TRUE   773    48   821
    Sum 10857    49 10906
```

DESeq

Simon Anders and Wolfgang Huber (2010): Differential expression analysis for sequence count data.Genome Biology 11:R106