

MIS-DESeq2

Sunit Jain

January 6, 2015

Contents

Dependencies	2
Generate a read count matrix using <code>htseq-count</code>	2
Merging duplicate genes	2
Import Counts into DESeq2	2
Exploring the Dataset	3
The <code>rlog</code> transformation	3
Sample distances	3
Poisson Distance	5
PCA plot	6
MDS plot	7
Counts	9
Raw vs Normalized Counts	9
Rank Abundance	13
Day Counts	13
Night Counts	14
Day vs Night	15
Differential Expression	16
Results	16
Multiple testing	16
Cook's Distance per gene	17
Diagnostic Plots	18
Plot Counts	18
MA-Plots	19
Dispersion Estimataion	20
P-Value Histogram	21
Gene clustering	22
Independent Filtering	24
Session Info	26

Dependencies

If you're unsure that you have all the packages required to run this workflow. Open the `Rmd` file in your favorite text editor (I used [RStudio](#)) and change the next line from `eval=FALSE` to `eval=TRUE`. Now, when you run this workflow, the dependencies should be installed first.

Generate a read count matrix using htseq-count

Sample command: `htseq-count -f bam -r name -t CDS -o scaffold.htseq.sam -i ID -q scaffold_sortedByName.b
all_combined.gff` This command was run for each sample individually.

Merging duplicate genes

I performed a self blast and looked at results that had a percent identity greater than 98%, query coverage greater than 96% and a minimum alignment length of 500 bases. Once I had this subset, I screened out the hits to exons since we won't be considering them for this experiment anyway. I was left with the following two gene pairs:

- scaffold_344578 MIS_1109813.1 scaffold_133898 MIS_10093600.14
- scaffold_219988 MIS_10179608.12 scaffold_555373 MIS_1172265.1

that had high enough similarity based on the thresholds mentioned above that their count data needed to be merged. The perl script `mergeCounts.pl` was run on each htseq-count output individually in order to accomplish this. Here is a sample command used for one of the htseq-count outputs: `perl mergeCounts.pl -l realDuplicateGenes.list -tsv Day_1.htseqCount.tsv -o Day_1.htseqCount.merged.tsv` where, `realDuplicateGenes.list` contains the two gene pairs mentioned above.

Import Counts into DESeq2

Once we were satisfied with the genes and their counts. We imported the count data into DESeq2.

Exploring the Dataset

The rlog transformation

Many common statistical methods for exploratory analysis of multidimensional data, especially methods for clustering and ordination (e.g., principal-component analysis and the like), work best for (at least approximately) homoskedastic data; this means that the variance of an observed quantity (here, the expression strength of a gene) does not depend on the mean. In RNA-Seq data, however, variance grows with the mean. For example, if one performs PCA (principal components analysis) directly on a matrix of normalized read counts, the result typically depends only on the few most strongly expressed genes because they show the largest absolute differences between samples. A simple and often used strategy to avoid this is to take the logarithm of the normalized count values plus a small pseudocount; however, now the genes with low counts tend to dominate the results because, due to the strong Poisson noise inherent to small count values, they show the strongest relative differences between samples.

As a solution, DESeq2 offers the regularized-logarithm transformation, or rlog for short. For genes with high counts, the rlog transformation differs not much from an ordinary log2 transformation. For genes with lower counts, however, the values are shrunken towards the genes' averages across all samples. Using an empirical Bayesian prior on inter-sample differences in the form of a ridge penalty, this is done such that the rlog-transformed data are approximately homoskedastic.

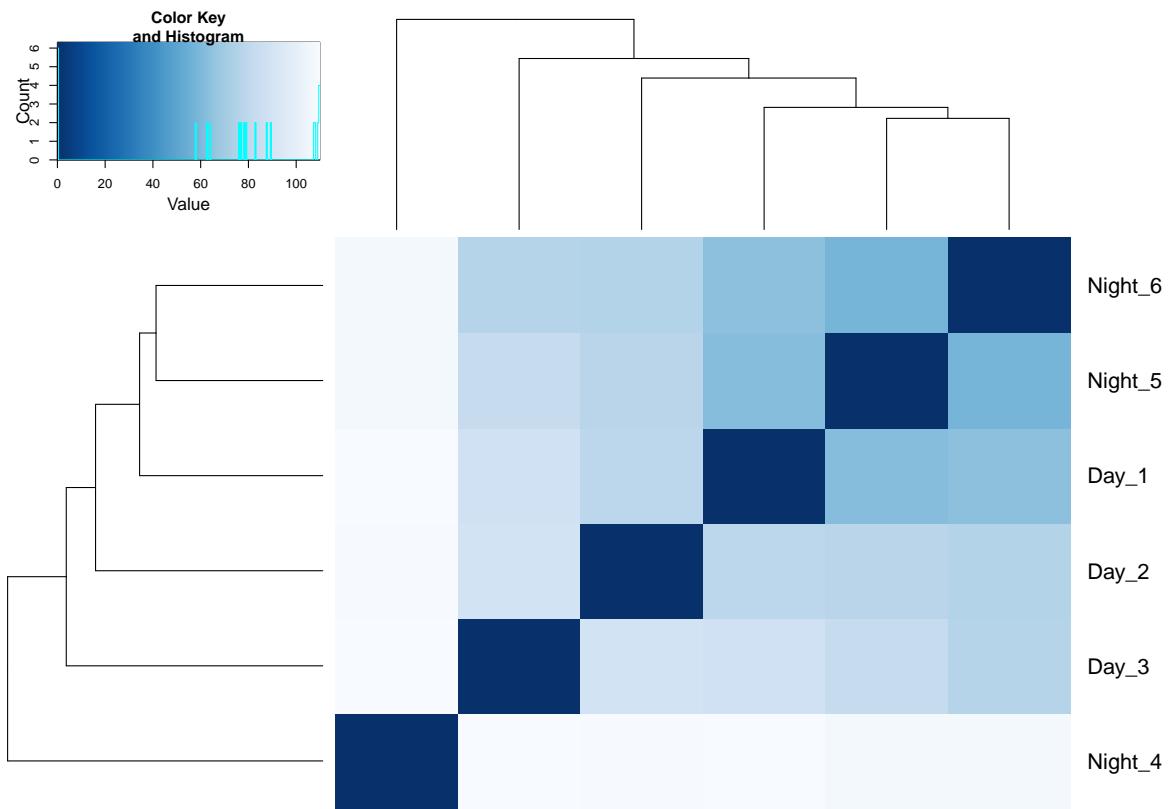
Note: the rlog transformation is provided for applications other than differential testing. For differential testing we recommend the DESeq function applied to raw counts, as described later in this workflow, which also takes into account the dependence of the variance of counts on the mean value during the dispersion estimation step.

Sample distances

A useful first step in an RNA-Seq analysis is often to assess overall similarity between samples: Which samples are similar to each other, which are different? Does this fit to the expectation from the experiment's design? We use the R function `dist` to calculate the Euclidean distance between samples. To avoid that the distance measure is dominated by a few highly variable genes, and have a roughly equal contribution from all genes, we use it on the rlog-transformed data:

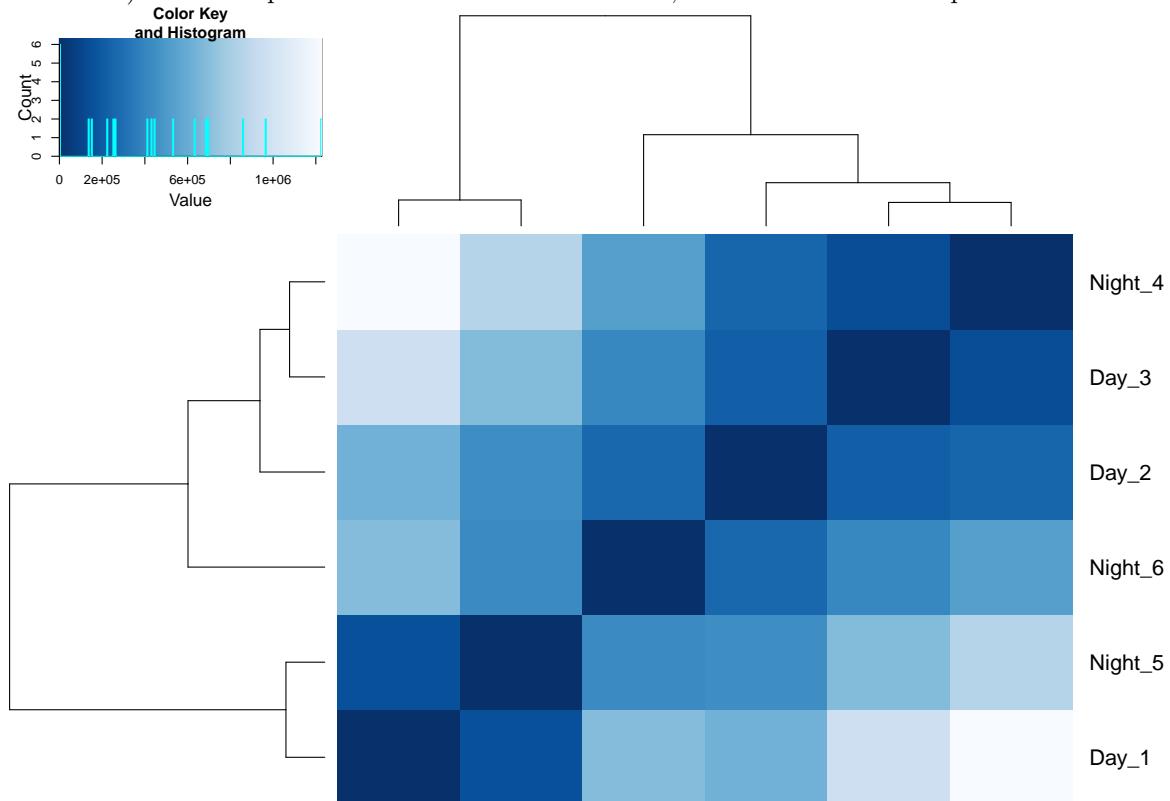
```
##          Day_1      Day_2      Day_3     Night_4     Night_5
## Day_2    79.20727
## Day_3    87.58113  89.42325
## Night_4 109.69210 109.28771 109.87333
## Night_5  62.61998  78.02932  82.92708 107.95831
## Night_6  63.86110  76.10091  76.87496 107.60440  58.12165
```

We visualize the distances in a heatmap:



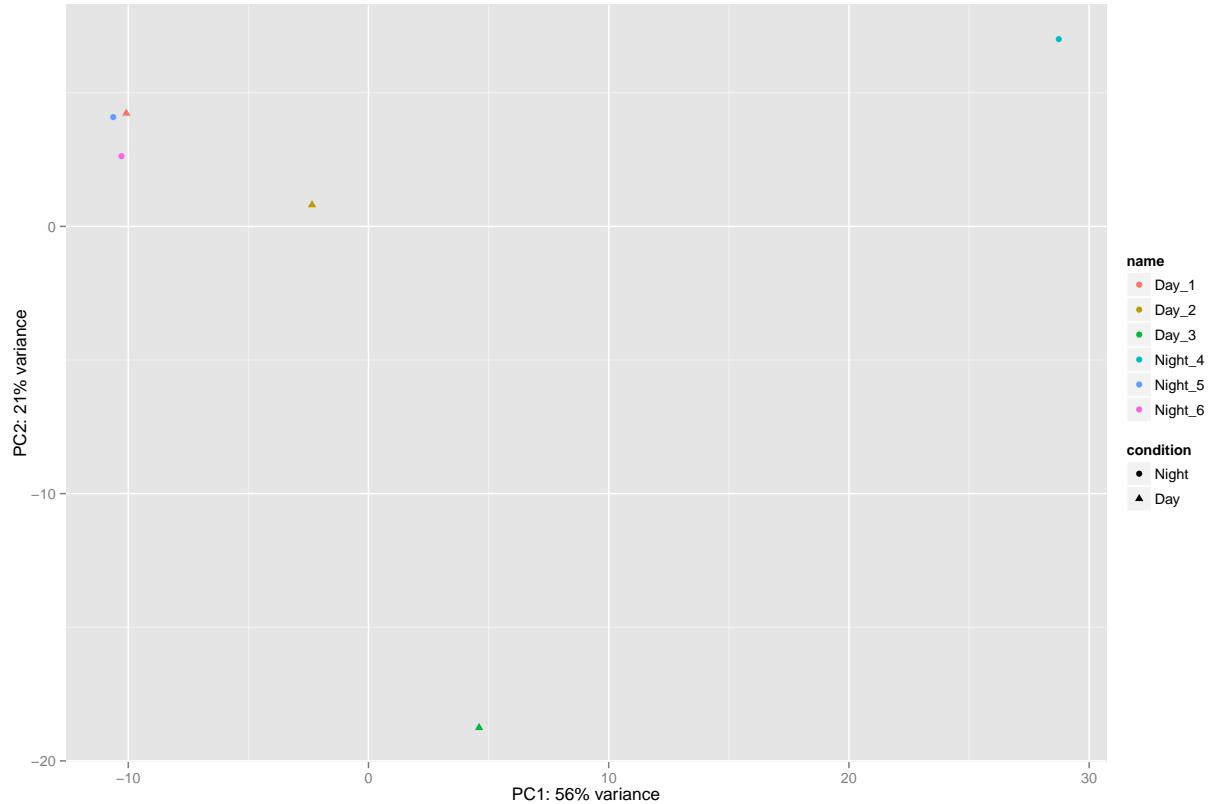
Poisson Distance

Another option for calculating sample distances is to use the Poisson Distance, implemented in the CRAN package PoiClaClu. Similar to the transformations offered in DESeq2, this measure of dissimilarity also takes the variance structure of counts into consideration when calculating the distances between samples. The PoissonDistance function takes the original count matrix (not normalized) with samples as rows instead of columns, so we need to transpose the counts in dds.



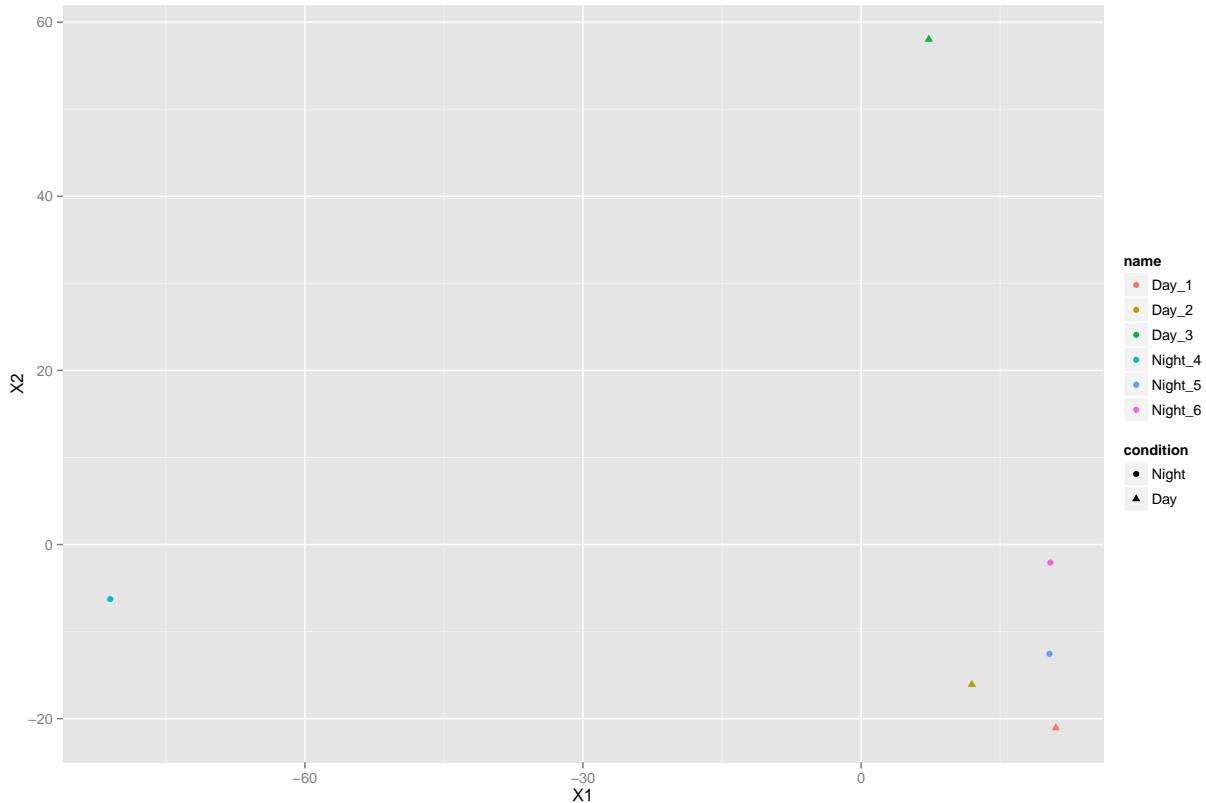
PCA plot

Another way to visualize sample-to-sample distances is a principal-components analysis (PCA). In this ordination method, the data points (i.e., here, the samples) are projected onto the 2D plane such that they spread out in the two directions which explain most of the differences in the data. The x-axis is the direction (or principal component) which separates the data points the most. The amount of the total variance which is contained in the direction is printed in the axis label.

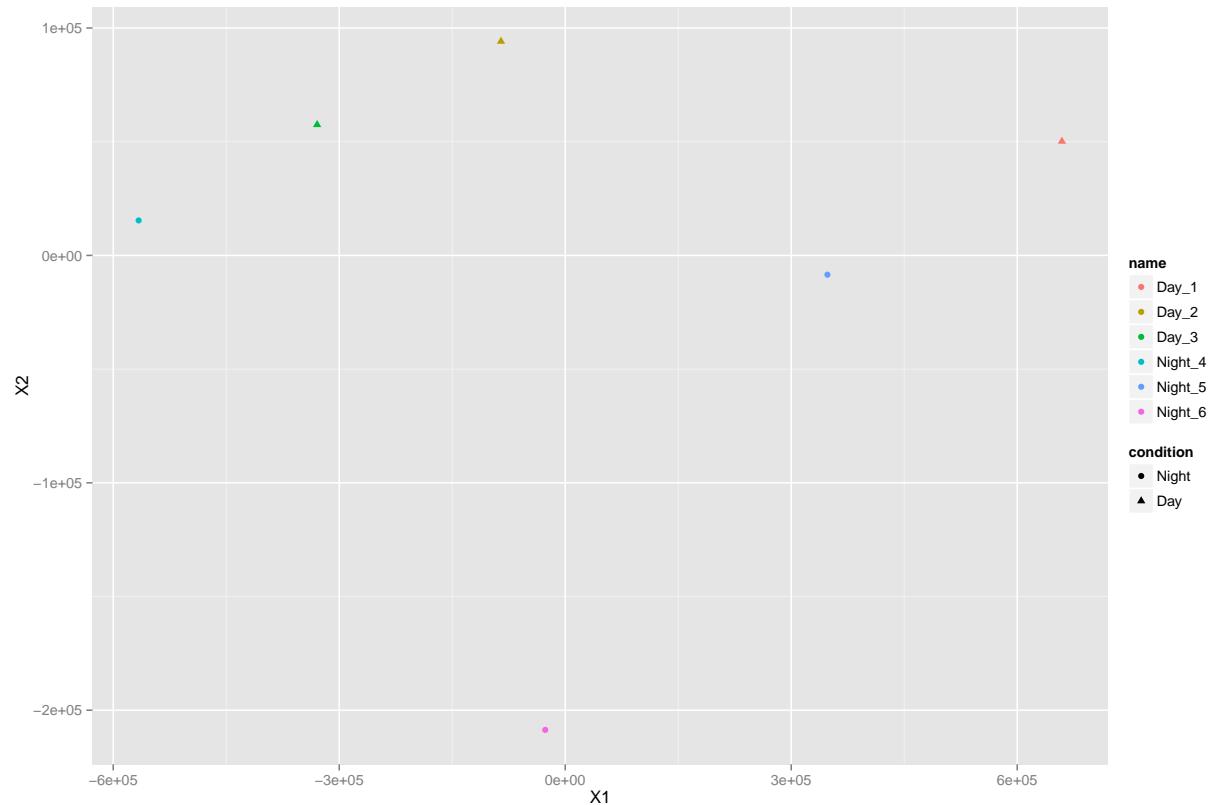


MDS plot

Another plot, very similar to the PCA plot, can be made using the multidimensional scaling (MDS) function in base R. This is useful when we don't have the original data, but only a matrix of distances. Here we have the MDS plot for the distances calculated from the rlog transformed counts:



And here from the PoissonDistance:



Counts

In order to normalise the raw counts we will start by determining the relative library sizes, or size factors for each library. For example, if the counts of the expressed genes in one sample are, on average, twice as high as in another, the size factor for the first sample should be twice as large as the one for the other sample. These size factors can be obtained with the function `estimateSizeFactors`:

```
##      Day_1      Day_2      Day_3     Night_4    Night_5    Night_6
## 1.6219205 0.6645299 0.5923596 0.3384058 2.2753012 2.5523719
```

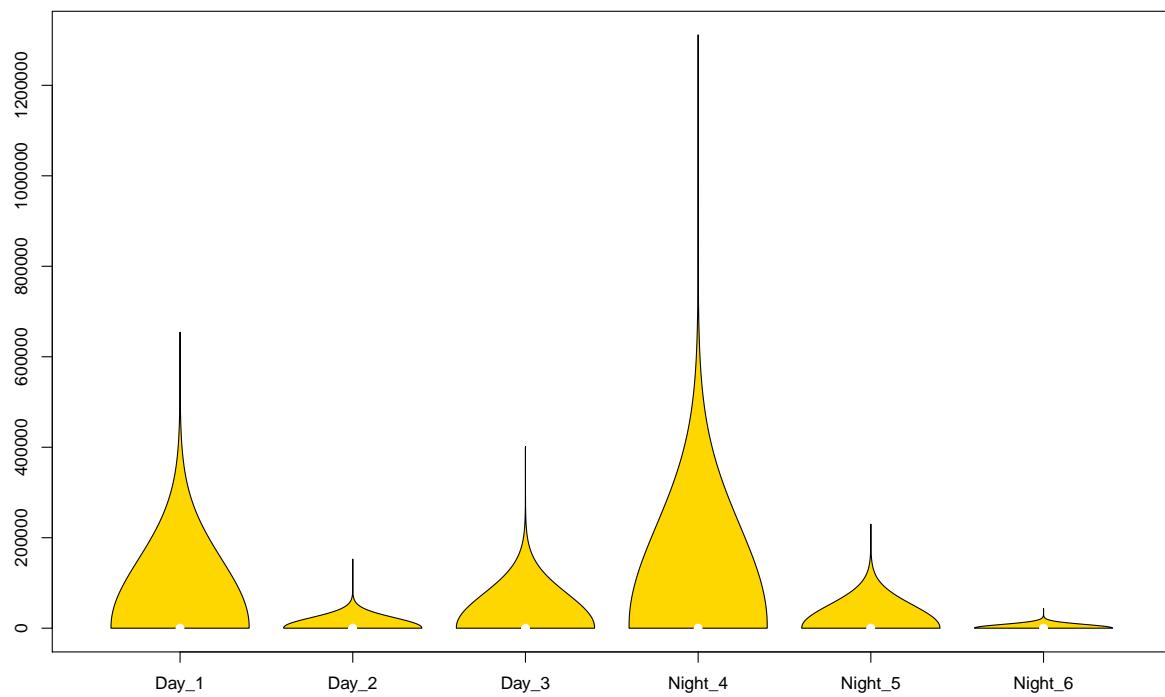
Raw vs Normalized Counts

Once we have this information, the normalised data is obtained by dividing each column of the count table by the corresponding size factor. We can perform this calculation by calling the function `counts` with a the `normalized` argument set as `TRUE`. Since we won't be normalizing this data, we'll set it as `FALSE`

Normalized

```
##                                     Day_1     Day_2 Day_3   Night_4   Night_5
## scaffold_0__MIS_10000001.1       0 1.504823      0 2.955032 0.0000000
## scaffold_0__MIS_10000001.1001    0 0.000000      0 0.000000 0.0000000
## scaffold_0__MIS_10000001.1005    0 0.000000      0 0.000000 0.0000000
## scaffold_0__MIS_10000001.103     0 0.000000      0 0.000000 0.0000000
## scaffold_0__MIS_10000001.107     0 0.000000      0 0.000000 0.4395022
## scaffold_0__MIS_10000001.111     0 0.000000      0 0.000000 0.0000000
##                                     Night_6
## scaffold_0__MIS_10000001.1      3.526132
## scaffold_0__MIS_10000001.1001    0.000000
## scaffold_0__MIS_10000001.1005    0.000000
## scaffold_0__MIS_10000001.103     0.000000
## scaffold_0__MIS_10000001.107     0.000000
## scaffold_0__MIS_10000001.111     0.000000
```

Violin Plots for Normalized Counts



```

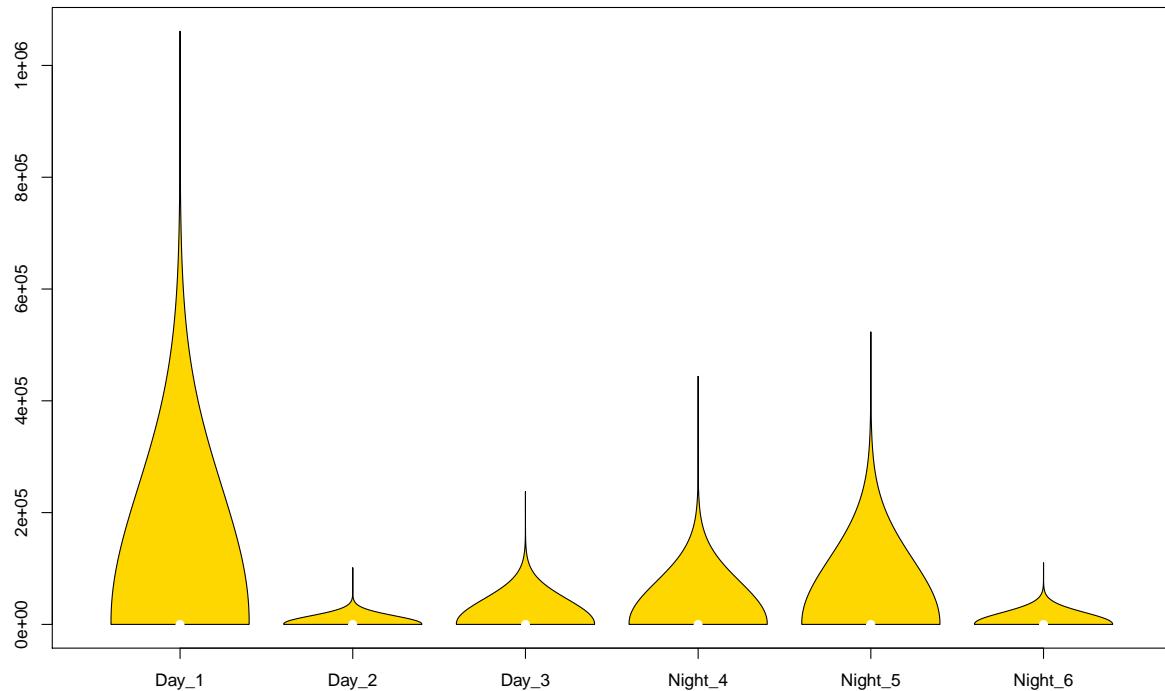
##      Day_1          Day_2          Day_3
## Min.   : 0.0   Min.   : 0.00   Min.   : 0.0
## 1st Qu.: 0.0   1st Qu.: 0.00   1st Qu.: 0.0
## Median : 0.0   Median : 0.00   Median : 0.0
## Mean   : 1.2   Mean   : 0.51   Mean   : 0.8
## 3rd Qu.: 0.0   3rd Qu.: 0.00   3rd Qu.: 0.0
## Max.   :654306.4 Max.   :152366.37 Max.   :400933.8
##      Night_4         Night_5        Night_6
## Min.   : 0.0   Min.   :0.0e+00   Min.   : 0.00
## 1st Qu.: 0.0   1st Qu.:0.0e+00   1st Qu.: 0.00
## Median : 0.0   Median :0.0e+00   Median : 0.00
## Mean   : 2.2   Mean   :4.8e-01   Mean   : 0.26
## 3rd Qu.: 0.0   3rd Qu.:0.0e+00   3rd Qu.: 0.00
## Max.   :1311227.4 Max.   :2.3e+05   Max.   :43083.85

```

Raw

```
##                                     Day_1 Day_2 Day_3 Night_4 Night_5 Night_6
## scaffold_0__MIS_10000001.1       0     1     0      1      0     9
## scaffold_0__MIS_10000001.1001    0     0     0      0      0     0
## scaffold_0__MIS_10000001.1005    0     0     0      0      0     0
## scaffold_0__MIS_10000001.103     0     0     0      0      0     0
## scaffold_0__MIS_10000001.107     0     0     0      0      1     0
## scaffold_0__MIS_10000001.111     0     0     0      0      0     0
```

Violin Plots for Raw Counts



```
##          Day_1           Day_2           Day_3
##  Min.   : 0.0   Min.   : 0.00   Min.   : 0.00
##  1st Qu.: 0.0   1st Qu.: 0.00   1st Qu.: 0.00
##  Median : 0.0   Median : 0.00   Median : 0.00
##  Mean   : 1.9   Mean   : 0.34   Mean   : 0.47
##  3rd Qu.: 0.0   3rd Qu.: 0.00   3rd Qu.: 0.00
##  Max.   :1061233.0  Max.   :101252.00  Max.   :237497.00
##          Night_4          Night_5          Night_6
##  Min.   : 0.0   Min.   : 0.0   Min.   :0.0e+00
##  1st Qu.: 0.0   1st Qu.: 0.0   1st Qu.:0.0e+00
##  Median : 0.0   Median : 0.0   Median :0.0e+00
##  Mean   : 0.8   Mean   : 1.1   Mean   :6.6e-01
##  3rd Qu.: 0.0   3rd Qu.: 0.0   3rd Qu.:0.0e+00
##  Max.   :443727.0  Max.   :523355.0  Max.   :1.1e+05
```

Highest Counts Here are the genes with the highest counts in each sample:

Day1

```
##                                     Day_1 Day_2 Day_3 Night_4 Night_5 Night_6
## scaffold_64771__MIS_10025904.6 1061233 57193 41606    97857  523355 109966
```

Day2

```
##                                     Day_1 Day_2 Day_3 Night_4 Night_5 Night_6
## scaffold_92040__MIS_10052095.6  3335 101252 41927   126590   4406   73889
```

Day3

```
##                                     Day_1 Day_2 Day_3 Night_4 Night_5 Night_6
## scaffold_345287__MIS_1110485.1 2299 15932 237497  443727   5927   24223
```

Night4

```
##                                     Day_1 Day_2 Day_3 Night_4 Night_5 Night_6
## scaffold_345287__MIS_1110485.1 2299 15932 237497  443727   5927   24223
```

Night5

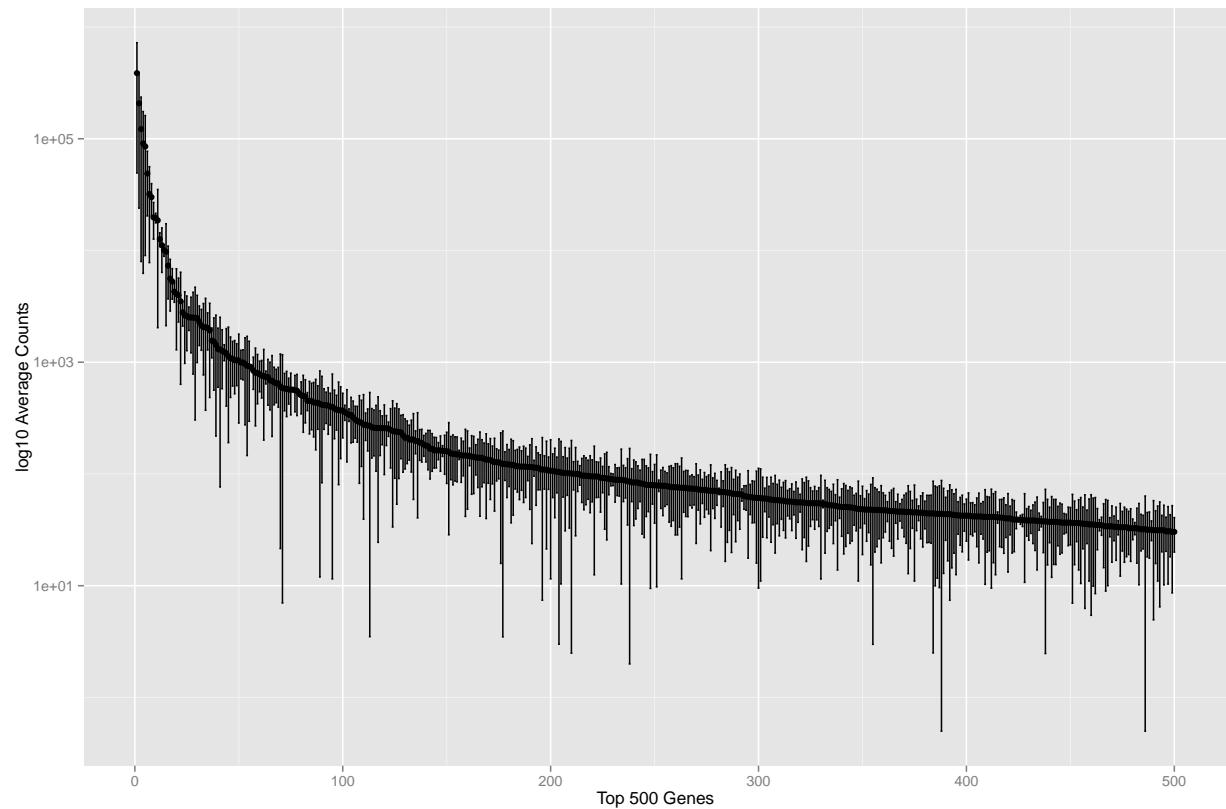
```
##                                     Day_1 Day_2 Day_3 Night_4 Night_5 Night_6
## scaffold_64771__MIS_10025904.6 1061233 57193 41606    97857  523355 109966
```

Night6

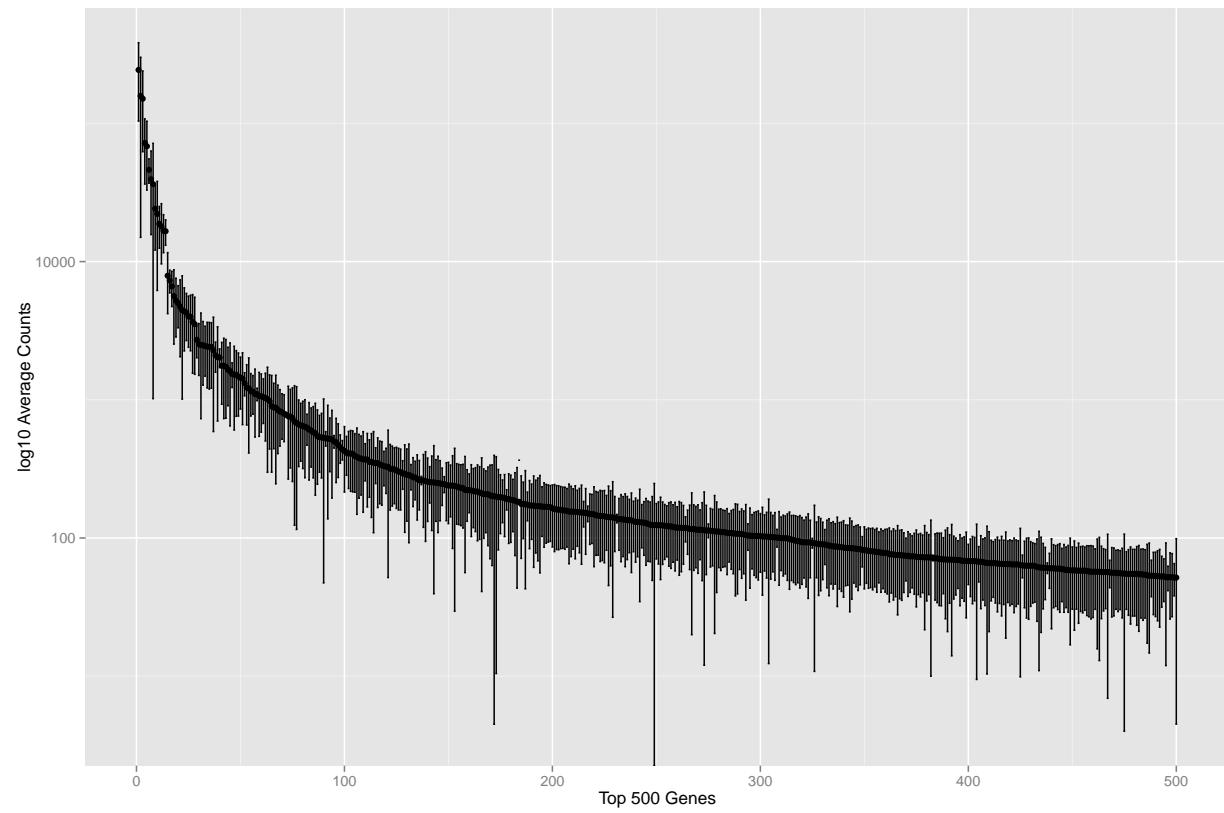
```
##                                     Day_1 Day_2 Day_3 Night_4 Night_5 Night_6
## scaffold_64771__MIS_10025904.6 1061233 57193 41606    97857  523355 109966
```

Rank Abundance

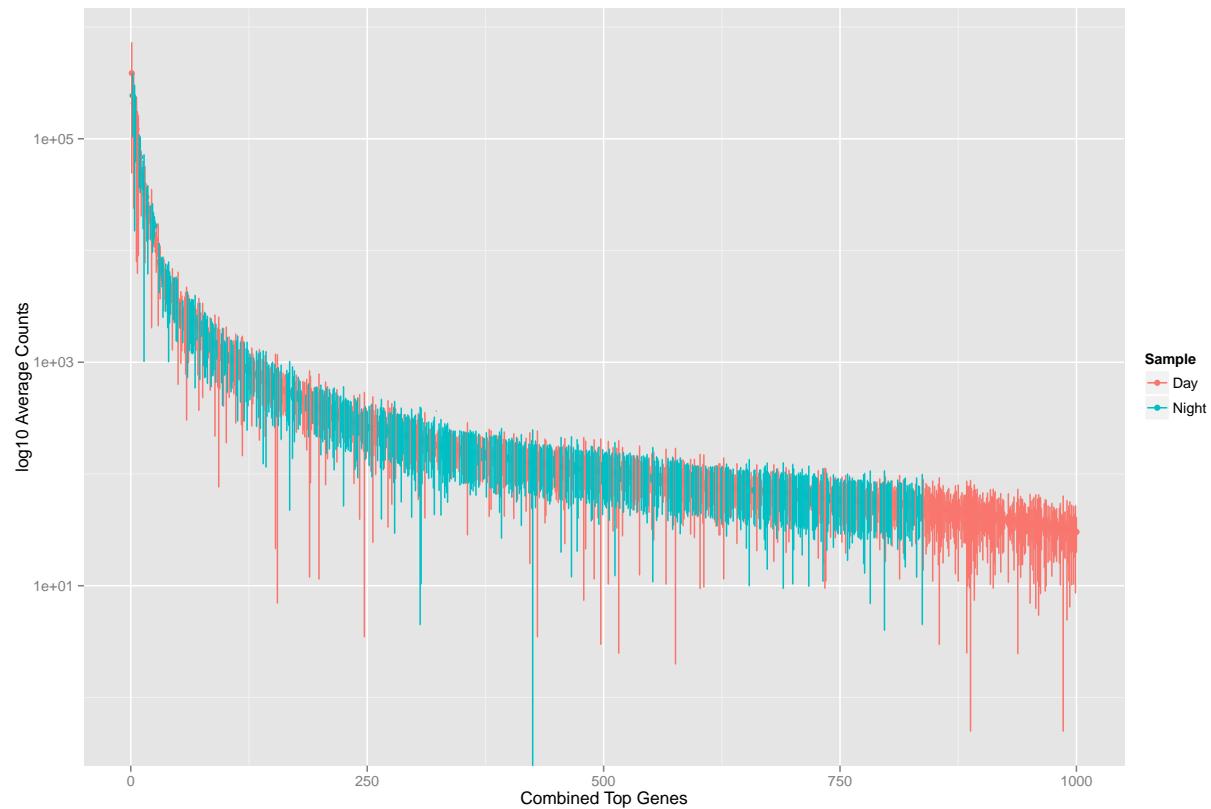
Day Counts



Night Counts



Day vs Night



Differential Expression

Differential expression was calculated using the DESeq2 wrapper function over 4 processors.

Results

As `res` is a DataFrame object, it carries metadata with information on the meaning of the columns:

```
## DataFrame with 6 rows and 2 columns
##           type                               description
## <character> <character>
## baseMean     intermediate      mean of normalized counts for all samples
## log2FoldChange results log2 fold change (MAP): condition Day vs Night
## lfcSE         results       standard error: condition Day vs Night
## stat          results      Wald statistic: condition Day vs Night
## pvalue        results      Wald test p-value: condition Day vs Night
## padj          results      BH adjusted p-values

##
## out of 86198 with nonzero total read count
## adjusted p-value < 0.1
## LFC > 0 (up)    : 0, 0%
## LFC < 0 (down)  : 0, 0%
## outliers [1]    : 95, 0.11%
## low counts [2]  : 0, 0%
## (mean count < 0)
## [1] see 'cooksCutoff' argument of ?results
## [2] see 'independentFiltering' argument of ?results
```

Multiple testing

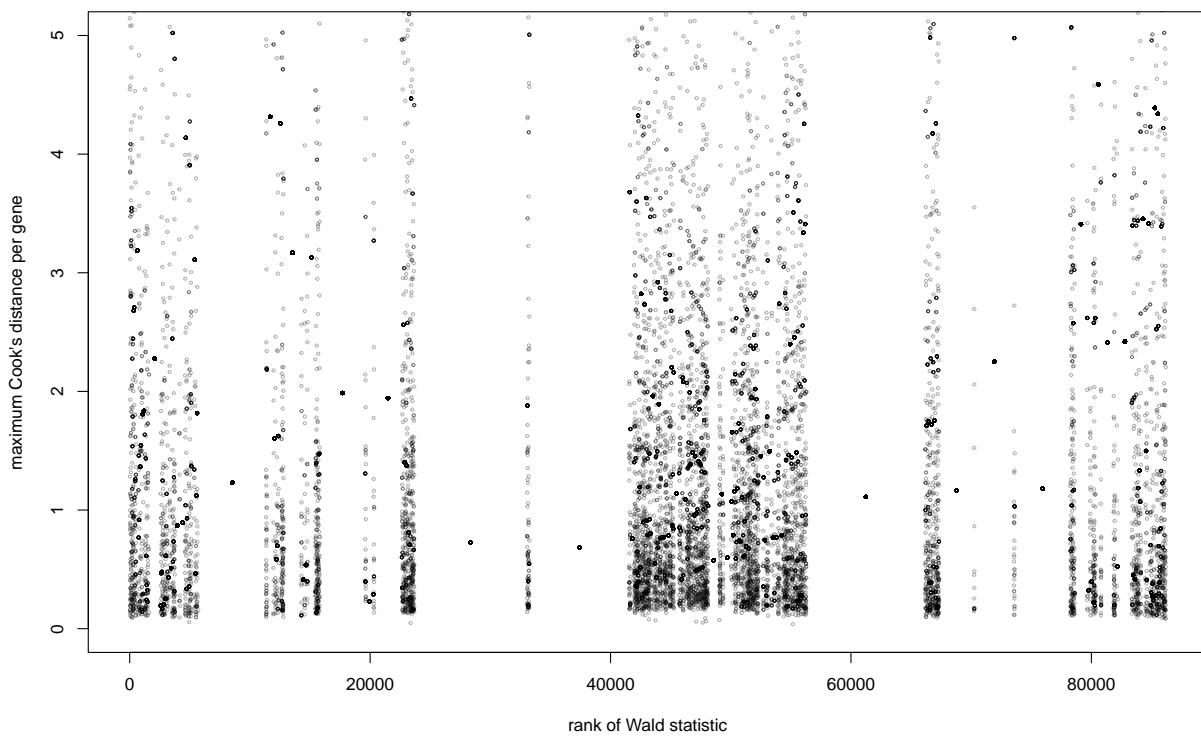
Novices in high-throughput biology often assume that thresholding these p values at a low value, say 0.05, as is often done in other settings, would be appropriate – but it is not. We briefly explain why: There are 224 genes with a p value below 0.05 among the 86103 genes, for which the test succeeded in reporting a p value.

DESeq2 uses the Benjamini-Hochberg (BH) adjustment as described in the base R `p.adjust` function; in brief, this method calculates for each gene an adjusted p value which answers the following question: if one called significant all genes with a p value less than or equal to this gene's p value threshold, what would be the fraction of false positives (the false discovery rate, FDR) among them (in the sense of the calculation outlined above)? These values, called the BH-adjusted p values, are given in the column `padj` of the `res` object. Hence, if we consider a fraction of 10% false positives acceptable, we can consider all genes with an adjusted p value below 10% = 0.1 as significant. How many such genes are there?

```
## [1] 0
```

Nothing! Infact, we get 0 genes with the adjusted p-value less than or equal to 0.79 and then suddenly at 0.8 we get 73523

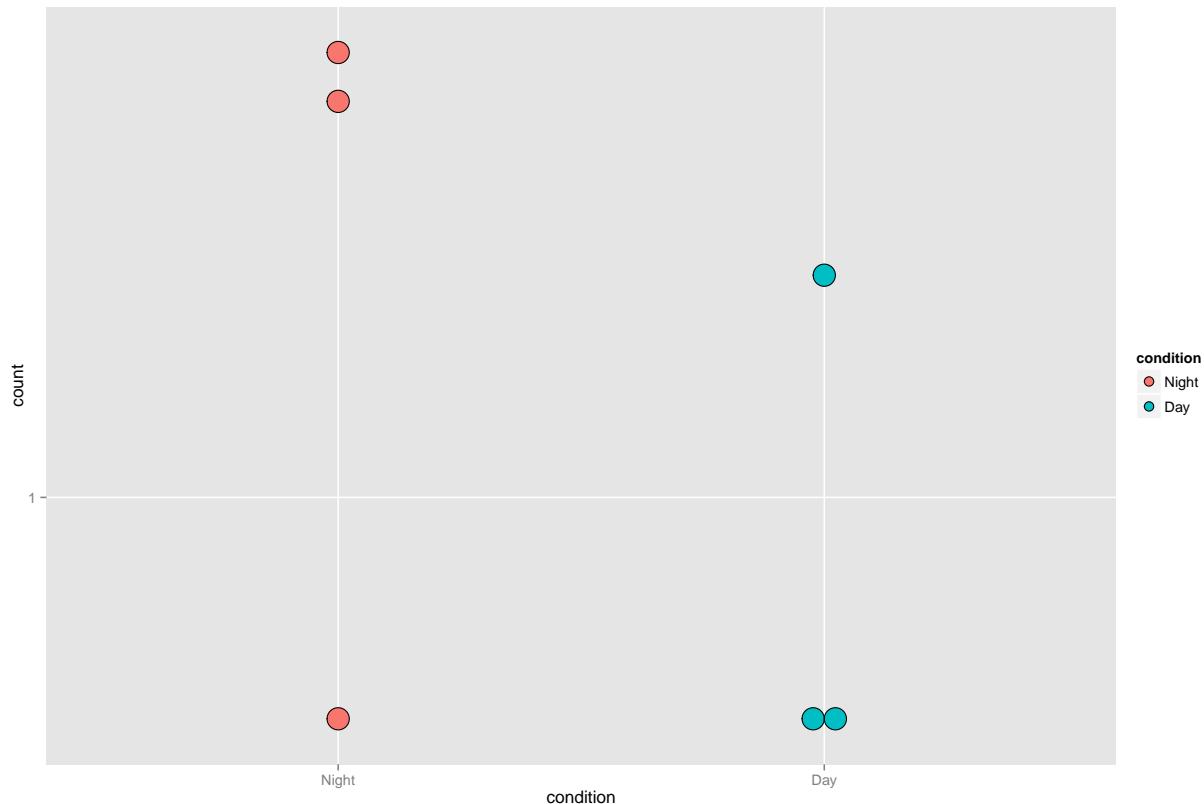
Cook's Distance per gene



Diagnostic Plots

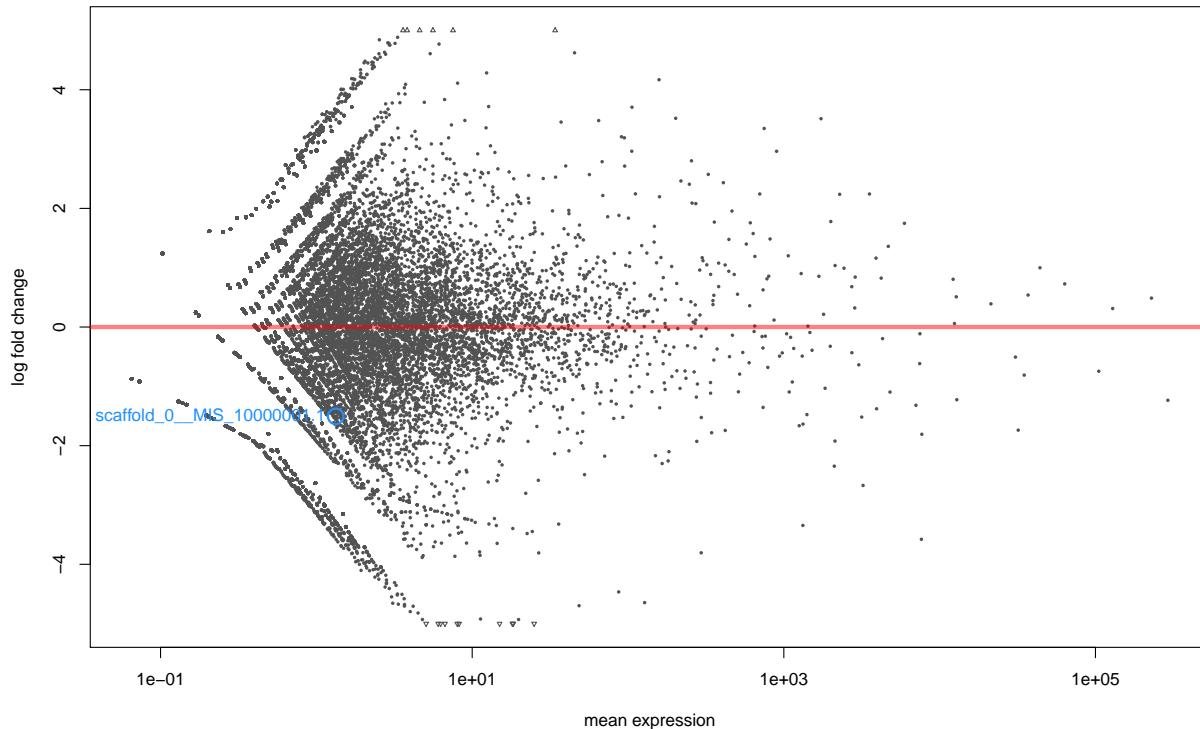
Plot Counts

A quick way to visualize the counts for a particular gene is to use the plotCounts function, which takes as arguments the DESeqDataSet, a gene name, and the group over which to plot the counts.



MA-Plots

An “MA-plot” provides a useful overview for an experiment with a two-group comparison. The log₂ fold change for a particular comparison is plotted on the y-axis and the average of the counts normalized by size factor is shown on the x-axis (“M” for minus, because a log ratio is equal to log minus log, and “A” for average).

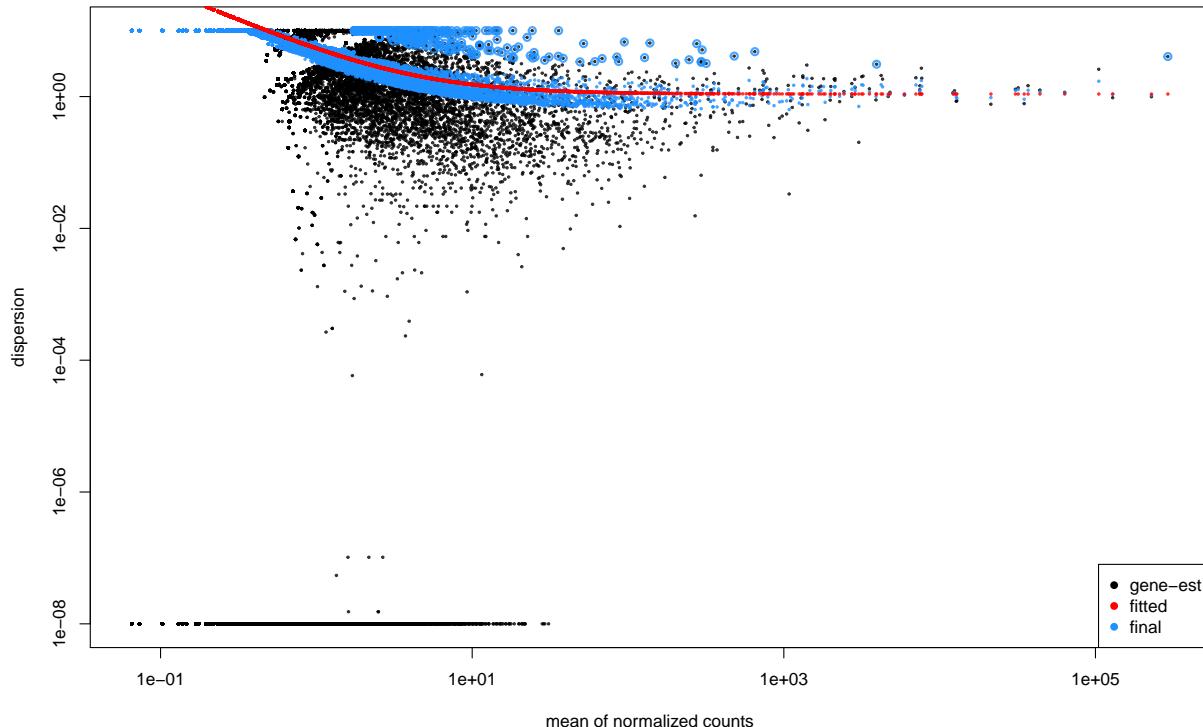


Each gene is represented with a dot. Genes with an adjusted p-value below a threshold (here 0.1, the default) are shown in red. The DESeq2 package incorporates a prior on log₂ fold changes, resulting in moderated log₂ fold changes from genes with low counts and highly variable counts, as can be seen by the narrowing of spread of points on the left side of the plot. This plot demonstrates that only genes with a large average normalized count contain sufficient information to yield a significant call.

Dispersion Estimation

Whether a gene is called significant depends not only on its LFC but also on its within-group variability, which DESeq2 quantifies as the dispersion. For strongly expressed genes, the dispersion can be understood as a squared coefficient of variation: a dispersion value of 0.01 means that the gene's expression tends to differ by typically $\text{sqrt}(0.01)=10\%$ between samples of the same treatment group. For weak genes, the Poisson noise is an additional source of noise.

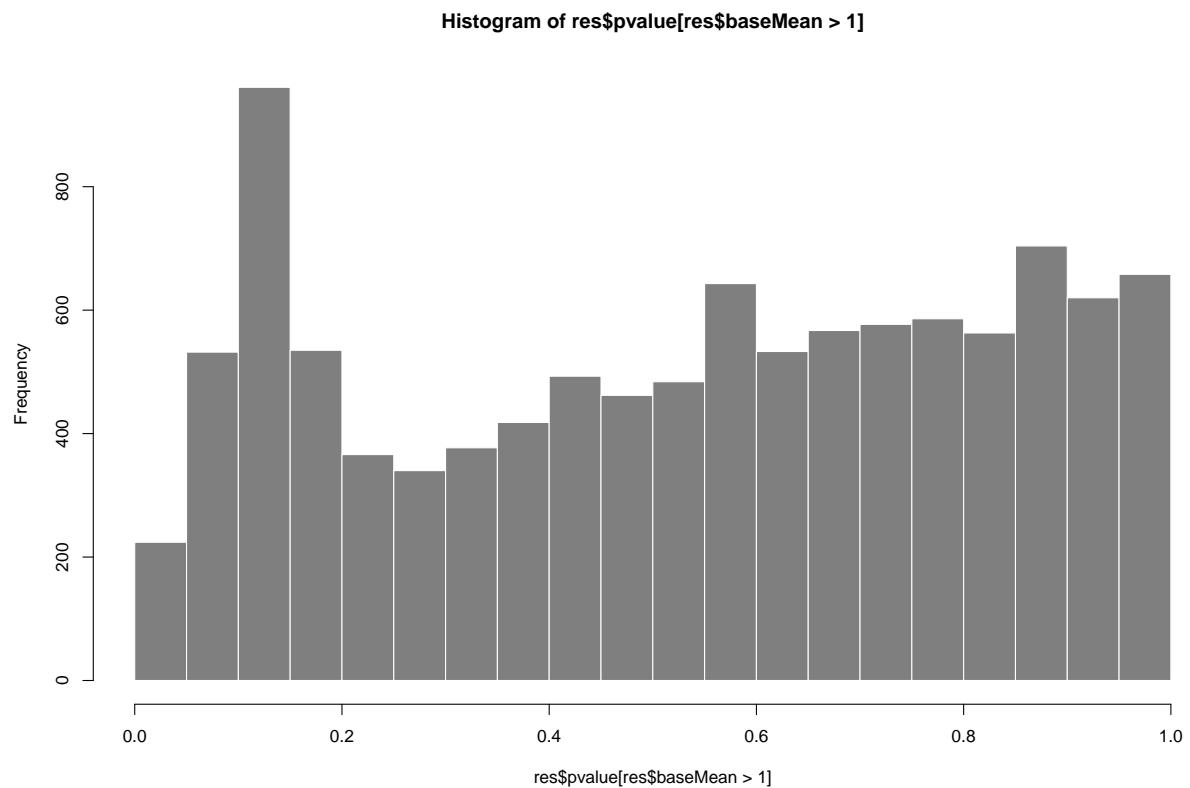
The function `plotDispEsts` visualizes DESeq2's dispersion estimates:



The black points are the dispersion estimates for each gene as obtained by considering the information from each gene separately. Unless one has many samples, these values fluctuate strongly around their true values. Therefore, we fit the red trend line, which shows the dispersions' dependence on the mean, and then shrink each gene's estimate towards the red line to obtain the final estimates (blue points) that are then used in the hypothesis test. The blue circles above the main “cloud” of points are genes which have high gene-wise dispersion estimates which are labelled as dispersion outliers. These estimates are therefore not shrunk toward the fitted trend line.

P-Value Histogram

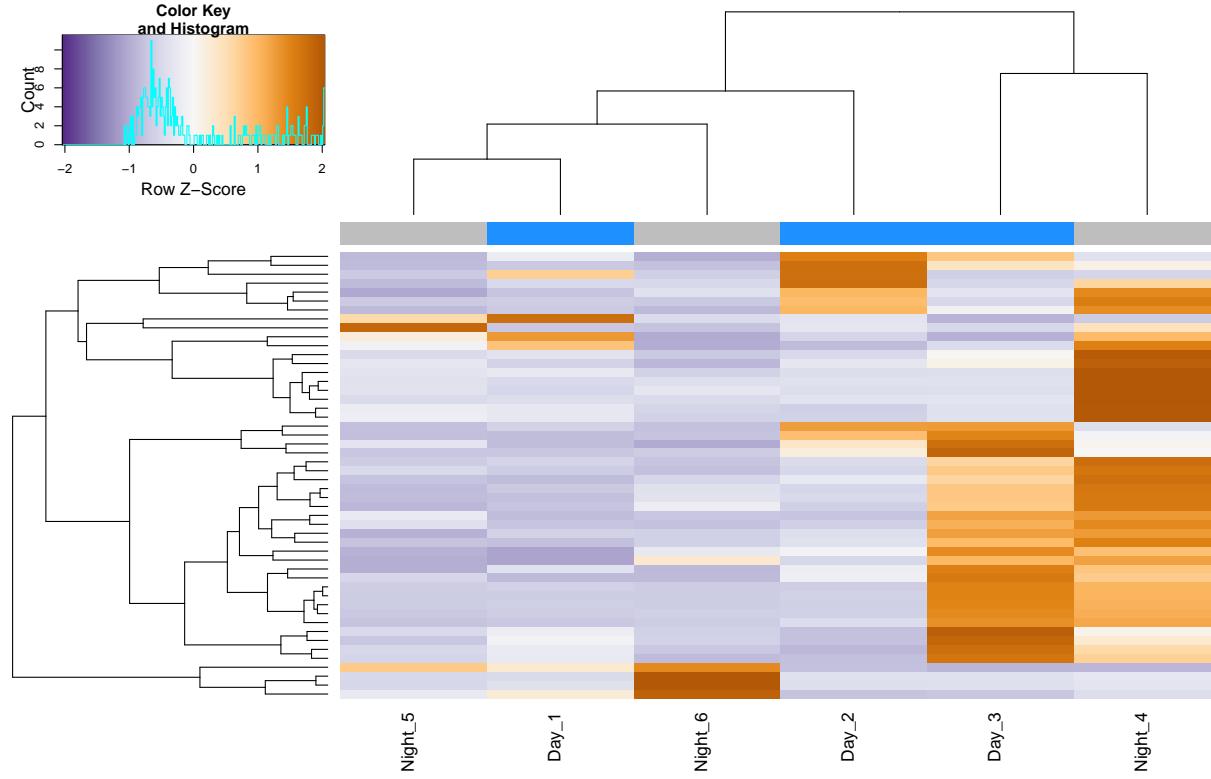
Another useful diagnostic plot is the histogram of the p values.



Gene clustering

In the sample distance heatmap made previously, the dendrogram at the side shows us a hierarchical clustering of the samples. Such a clustering can also be performed for the genes. Since the clustering is only relevant for genes that actually carry signal, one usually carries it out only for a subset of most highly variable genes. Here, for demonstration, let us select the 50 genes with the highest variance across samples. We will work with the rlog transformed counts:

The heatmap becomes more interesting if we do not look at absolute expression strength but rather at the amount by which each gene deviates in a specific sample from the gene's average across all samples. Hence, we center each genes' values across samples, and plot a heatmap. We provide the column side colors to help identify the treated samples (in blue) from the untreated samples (in grey).



We can now see blocks of genes which covary across patients. Note that a set of genes at the top of the heatmap are separating the N061011 cell line from the others. At the bottom of the heatmap, we see a set of genes for which the treated samples have higher gene expression.

```
## log2 fold change (MAP): condition Day vs Night
## Wald test p-value: condition Day vs Night
## DataFrame with 6 rows and 7 columns
##           baseMean log2FoldChange      lfcSE
## <numeric>     <numeric>     <numeric>
## scaffold_42417__MIS_10006030.1 34.069291    5.908385  1.477775
## scaffold_12010__MIS_10012011.1 45.377524    4.622566  1.236734
## scaffold_12010__MIS_10012011.12 158.180835   4.168822  1.166009
## scaffold_3287__MIS_10003288.28 14.991303   -6.269444  1.769826
## scaffold_200891__MIS_10160517.8  7.531982    6.118912  1.740125
## scaffold_12010__MIS_10012011.10 5.599365    5.841213  1.746641
```

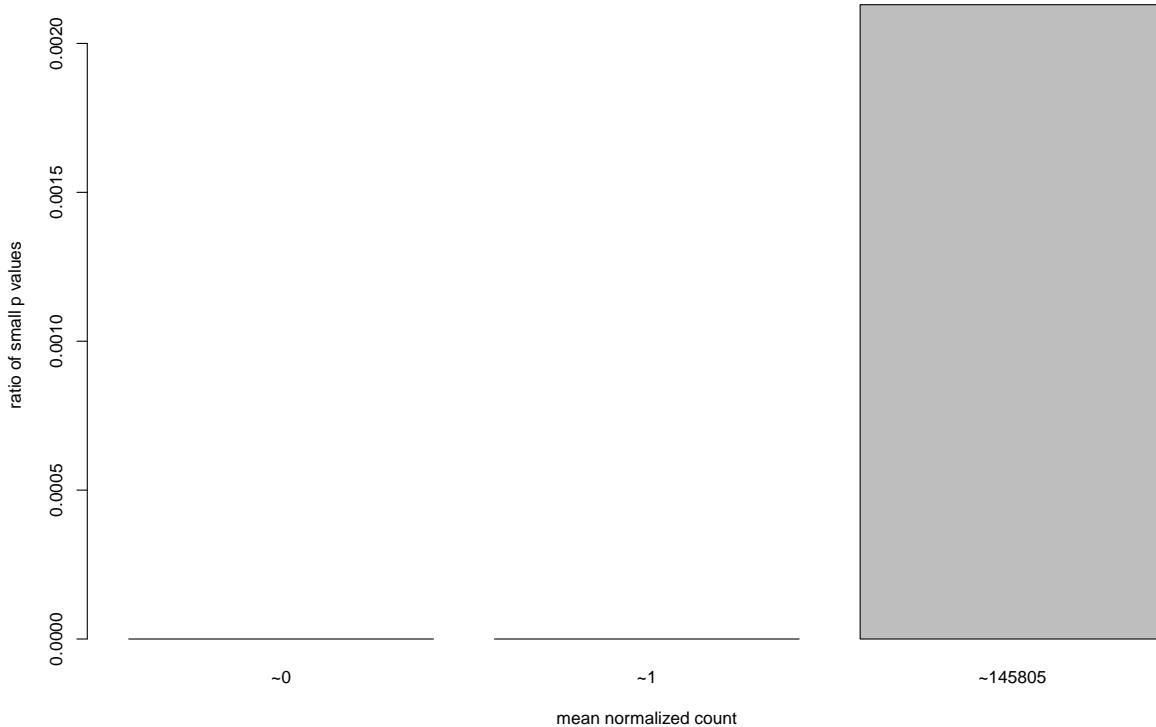
```

##                      stat      pvalue      padj
## <numeric>    <numeric>    <numeric>
## scaffold_42417__MIS_10006030.1 3.998163 6.383606e-05 0.793278
## scaffold_12010__MIS_10012011.1 3.737719 1.856971e-04 0.793278
## scaffold_12010__MIS_10012011.12 3.575293 3.498354e-04 0.793278
## scaffold_3287__MIS_10003288.28 -3.542406 3.964943e-04 0.793278
## scaffold_200891__MIS_10160517.8 3.516364 4.375005e-04 0.793278
## scaffold_12010__MIS_10012011.10 3.344255 8.250377e-04 0.793278
##                                     Name
## <character>
## scaffold_42417__MIS_10006030.1 scaffold_42417__MIS_10006030.1
## scaffold_12010__MIS_10012011.1 scaffold_12010__MIS_10012011.1
## scaffold_12010__MIS_10012011.12 scaffold_12010__MIS_10012011.12
## scaffold_3287__MIS_10003288.28 scaffold_3287__MIS_10003288.28
## scaffold_200891__MIS_10160517.8 scaffold_200891__MIS_10160517.8
## scaffold_12010__MIS_10012011.10 scaffold_12010__MIS_10012011.10

```

Independent Filtering

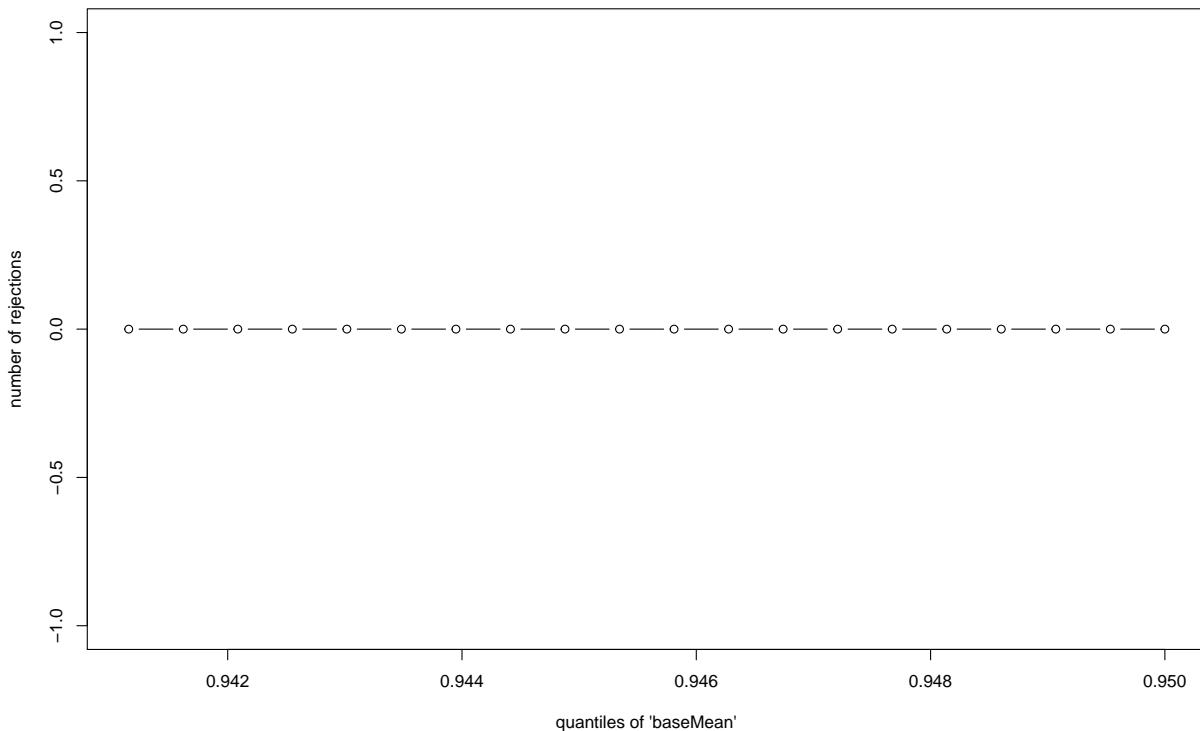
The MA plot highlights an important property of RNA-Seq data. For weakly expressed genes, we have no chance of seeing differential expression, because the low read counts suffer from so high Poisson noise that any biological effect is drowned in the uncertainties from the read counting. We can also show this by examining the ratio of small p values (say, less than, 0.01) for genes binned by mean normalized count:



At first sight, there may seem to be little benefit in filtering out these genes. After all, the test found them to be non-significant anyway. However, these genes have an influence on the multiple testing adjustment, whose performance improves if such genes are removed. By removing the weakly-expressed genes from the input to the FDR procedure, we can find more genes to be significant among those which we keep, and so improved the power of our test. This approach is known as independent filtering.

The DESeq2 software automatically performs independent filtering which maximizes the number of genes which will have adjusted p value less than a critical value (by default, alpha is set to 0.1). This automatic independent filtering is performed by, and can be controlled by, the results function. We can observe how the number of rejections changes for various cutoffs based on mean normalized count. The following optimal threshold and table of possible values is stored as an attribute of the results object.

```
##      94.1155%
## 0.003842503
```



The term independent highlights an important caveat. Such filtering is permissible only if the filter criterion is independent of the actual test statistic. Otherwise, the filtering would invalidate the test and consequently the assumptions of the BH procedure. This is why we filtered on the average over all samples: this filter is blind to the assignment of samples to the treatment and control group and hence independent. The independent filtering software used inside DESeq2 comes from the genefilter package, which contains a reference to a paper describing the statistical foundation for independent filtering.

Session Info

```
## R version 3.2.0 (2015-04-16)
## Platform: x86_64-apple-darwin13.4.0 (64-bit)
## Running under: OS X 10.10.3 (Yosemite)
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## attached base packages:
## [1] parallel   stats4     stats      graphics   grDevices  utils      datasets
## [8] methods    base
##
## other attached packages:
## [1] genefilter_1.50.0      BiocParallel_1.2.0
## [3] dplyr_0.4.1           tidyR_0.2.0
## [5] vioplot_0.2            sm_2.2-5.4
## [7] ggplot2_1.0.1          PoiClaClu_1.0.2
## [9] RColorBrewer_1.1-2     gplots_2.16.0
## [11] ReportingTools_2.8.0   RSQLite_1.0.0
## [13] DBI_0.3.1              knitr_1.10
## [15] DESeq2_1.8.0           RcppArmadillo_0.5.000.0
## [17] Rcpp_0.11.5             GenomicRanges_1.20.3
## [19] GenomeInfoDb_1.4.0     IRanges_2.2.1
## [21] S4Vectors_0.6.0         BiocGenerics_0.14.0
##
## loaded via a namespace (and not attached):
## [1] Category_2.34.0          bitops_1.0-6
## [3] tools_3.2.0               rpart_4.1-9
## [5] KernSmooth_2.23-14       lazyeval_0.1.10
## [7] Hmisc_3.15-0              colorspace_1.2-6
## [9] nnet_7.3-9                gridExtra_0.9.1
## [11] GGally_0.5.0              graph_1.46.0
## [13] Biobase_2.28.0           formatR_1.2
## [15] labeling_0.3              rtracklayer_1.28.2
## [17] ggbio_1.16.0              caTools_1.17.1
## [19] scales_0.2.4              RBGL_1.44.0
## [21] stringr_0.6.2             digest_0.6.8
## [23] Rsamtools_1.20.1          foreign_0.8-63
## [25] rmarkdown_0.5.1            R.utils_2.0.0
## [27] AnnotationForge_1.10.0    XVector_0.8.0
## [29] dichromat_2.0-0           htmltools_0.2.6
## [31] limma_3.24.1              BSgenome_1.36.0
## [33] PFAM.db_3.1.2             GOstats_2.34.0
## [35] hwriter_1.3.2              gtools_3.4.2
## [37] acepack_1.3-3.3           R.oo_1.19.0
## [39] magrittr_1.5                VariantAnnotation_1.14.0
## [41] RCurl_1.95-4.6             GO.db_3.1.2
## [43] Formula_1.2-1              futile.logger_1.4.1
## [45] Matrix_1.2-0              munsell_0.4.2
## [47] proto_0.3-10              R.methodsS3_1.7.0
## [49] yaml_2.1.13                edgeR_3.10.0
## [51] MASS_7.3-40                 zlibbioc_1.14.0
## [53] plyr_1.8.2                  grid_3.2.0
```

```
## [55] gdata_2.13.3           lattice_0.20-31
## [57] Biostrings_2.36.0        splines_3.2.0
## [59] GenomicFeatures_1.20.0   annotate_1.46.0
## [61] locfit_1.5-9.1          geneplotter_1.46.0
## [63] reshape2_1.4.1          biomaRt_2.24.0
## [65] futile.options_1.0.0    XML_3.98-1.1
## [67] evaluate_0.7            biovizBase_1.16.0
## [69] latticeExtra_0.6-26     lambda.r_1.1.7
## [71] gtable_0.1.2            assertthat_0.1
## [73] reshape_0.8.5           xtable_1.7-4
## [75] survival_2.38-1         snow_0.3-13
## [77] OrganismDbi_1.10.0      GenomicAlignments_1.4.0
## [79] AnnotationDbi_1.30.0    cluster_2.0.1
## [81] GSEABase_1.30.1
```