

Statistical report of project TSSC_nam7_xrn1: pairwise comparison(s) of conditions with DESeq2

Author: Christophe Malabat

Date: 2015-03-30

The SARTools R package which generated this report has been developed at PF2 - Institut Pasteur by M.-A. Dillies and H. Varet (hugo.varet@pasteur.fr). Thanks to cite H. Varet, J.-Y. Coppee and M.-A. Dillies, *SARTools: a DESeq2- and edgeR-based R pipeline for comprehensive differential analysis of RNA-seq data*, 2014 (submitted) when using this tool for any analysis published.

Table of contents

1. Introduction
2. Description of raw data
3. Variability within the experiment: data exploration
4. Normalization
5. Differential analysis
6. R session information and parameters
7. Bibliography

1 Introduction

The analyses reported in this document are part of the TSSC_nam7_xrn1 project. The aim is to find features that are differentially expressed between WT, upf1, upf1_xrn1 and xrn1. The statistical analysis process includes data normalization, graphical exploration of raw and normalized data, test for differential expression for each feature between the conditions, raw p-value adjustment and export of lists of features having a significant differential expression between the conditions.

The analysis is performed using the R software [R Core Team, 2014], Bioconductor [Gentleman, 2004] packages including DESeq2 [Anders, 2010 and Love, 2014] and the SARTools package developed at PF2 - Institut Pasteur. Normalization and differential analysis are carried out according to the DESeq2 model and package. This report comes with additional tab-delimited text files that contain lists of differentially expressed features.

For more details about the DESeq2 methodology, please refer to its related publications [Anders, 2010 and Love, 2014].

2 Description of raw data

The count data files and associated biological conditions are listed in the following table.

label	files	group
LT_01.WT	LT_01.WT_for_desq2.txt	WT
LT_02.WT	LT_02.WT_for_desq2.txt	WT
LT_03.WT	LT_03.WT_for_desq2.txt	WT
LT_01.upf1	LT_01.upf1_for_desq2.txt	upf1
LT_02.upf1	LT_02.upf1_for_desq2.txt	upf1
LT_03.upf1	LT_03.upf1_for_desq2.txt	upf1
LT_01.upf1.xrn1	LT_01.upf1.xrn1_for_desq2.txt	upf1_xrn1
LT_02.upf1.xrn1	LT_02.upf1.xrn1_for_desq2.txt	upf1_xrn1
LT_03.upf1.xrn1	LT_03.upf1.xrn1_for_desq2.txt	upf1_xrn1
LT_01.xrn1	LT_01.xrn1_for_desq2.txt	xrn1
LT_02.xrn1	LT_02.xrn1_for_desq2.txt	xrn1
LT_03.xrn1	LT_03.xrn1_for_desq2.txt	xrn1

Table 1: Data files and associated biological conditions.

After loading the data we first have a look at the raw data table itself. The data table contains one row per annotated feature and one column per sequenced sample. Row names of this table are feature IDs (unique identifiers). The table contains raw count values representing the number of reads that map onto the features. For this project, there are 11353 features in the count data table.

	LT_01.WT	LT_02.WT	LT_03.WT	LT_01.upf1	LT_02.upf1	LT_03.upf1	LT_01.upf1.xrn1	LT_02.upf1.xrn1	LT_03.upf1.xrn1	LT_01.xrn1	LT_02.xrn1	LT_03.xrn1
15S_rRNA	22936	40732	63667	25754	57946	83700	24961	68353	93314			
21S_rRNA	141198	120244	261442	147422	168563	315985	114999	190164	305163			
CUT_GIM_0001	185	23	208	156	100	256	310	94	405			
CUT_GIM_0002	70	0	70	191	0	191	53	0	53			
CUT_GIM_0003	138	0	138	135	0	135	98	10	108			
CUT_GIM_0004	69	3	72	86	22	109	49	7	56			

Table 2: Partial view of the count data table.

Looking at the summary of the count table provides a basic description of these raw counts (min and max values, median, etc).

	LT_01.WT	LT_02.WT	LT_03.WT	LT_01.upf1	LT_02.upf1	LT_03.upf1	LT_01.upf1.xrn1	LT_02.upf1.xrn1	LT_03.upf1.xrn1	LT_01.xr
Min.	0	0	0	0	0	0	0	0	0	0
1st Qu.	200	77	297	343	219	575	330	223	562	600
Median	1401	1117	2590	2513	2650	5322	2708	2754	5531	2100
Mean	30586	32344	62930	29236	35103	64338	24969	32661	57630	33000
3rd Qu.	7647	6648	14463	7877	9581	17727	10149	11020	21103	13000
Max.	65705285	72691948	138397233	64340619	76372314	140712933	48937808	68329214	117267022	67006000

Table 3: Summary of the raw counts.

Figure 1 shows the total number of mapped reads for each sample. Reads that map on multiple locations on the transcriptome are counted more than once, as far as they are mapped on less than 50 different loci. We expect total read counts to be similar within conditions, they may be different across conditions. Total counts sometimes vary widely between replicates. This may happen for several reasons, including:

- different rRNA contamination levels between samples (even between biological replicates);
- slight differences between library concentrations, since they may be difficult to measure with high precision.

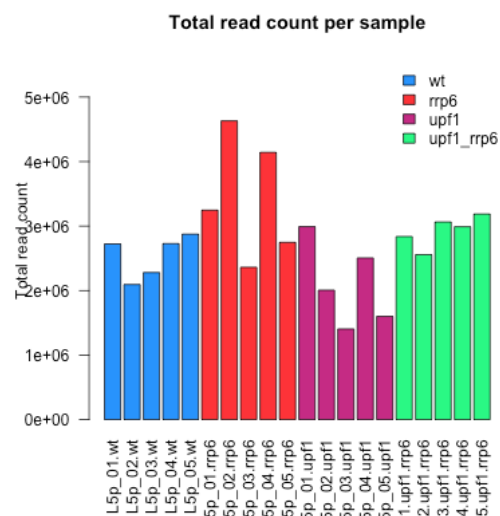


Figure 1: Number of mapped reads per sample. Colors refer to the biological condition of the sample.

Figure 2 shows the proportion of features with no read count in each sample. We expect this proportion to be similar within conditions. Features with null read counts in the 12 samples are left in the data but are not taken into account for the analysis with DESeq2. Here, 3 features (0.03%) are in this situation (dashed line). Results for those features (fold-change and p-values) are set to NA in the results files.

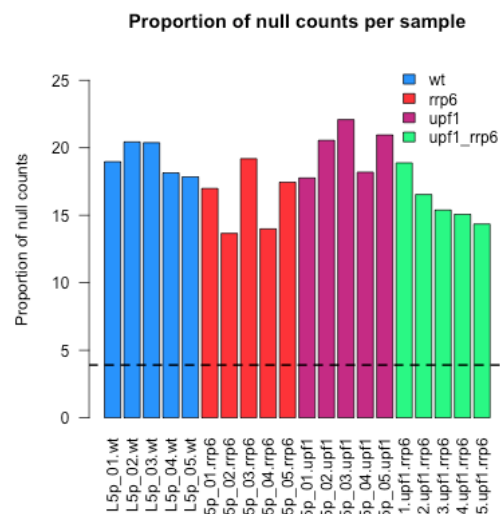


Figure 2: Proportion of features with null read counts in each sample.

Figure 3 shows the distribution of read counts for each sample. For sake of readability, $\log_2(\text{counts} + 1)$ are used instead of raw counts. Again we expect replicates to have similar distributions. In addition, this figure shows if read counts are preferably low, medium or high. This depends on the organisms as well as the biological conditions under consideration.

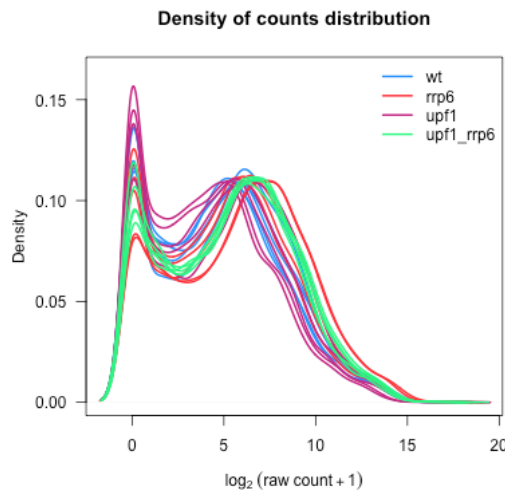


Figure 3: Density distribution of read counts.

It may happen that one or a few features capture a high proportion of reads (up to 20% or more). This phenomenon should not influence the normalization process. The DESeq2 normalization has proved to be robust to this situation [Dillies, 2012]. Anyway, we expect these high count features to be the same across replicates. They are not necessarily the same across conditions. Figure 4 and table 4 illustrate the possible presence of such high count features in the data set.

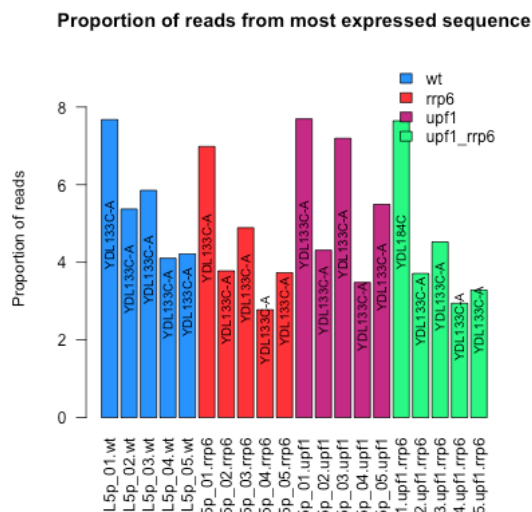


Figure 4: Percentage of reads associated with the sequence having the highest count (provided in each box on the graph) for each sample.

	LT_01.WT	LT_02.WT	LT_03.WT	LT_01.upf1	LT_02.upf1	LT_03.upf1	LT_01.upf1.xrn1	LT_02.upf1.xrn1	LT_03.upf1.xrn1	LT_01.xi
RDN37-2	18.92	19.80	19.37	19.39	19.16	19.26	17.26	18.43	17.92	17
RDN25-2	14.74	13.74	14.23	15.89	13.39	14.53	14.36	13.21	13.71	14
RDN37-1	11.41	12.43	11.94	11.48	11.97	11.75	10.33	11.43	10.95	10

Table 4: Percentage of reads associated with the sequences having the highest counts.

We may wish to assess the similarity between samples across conditions. A pairwise scatter plot is produced (figure 5) to show how replicates and samples from different biological conditions are similar or different ($\log_2(\text{counts} + 1)$ are used instead of raw count values). Moreover, as the Pearson correlation has been shown not to be relevant to measure the similarity between replicates, the SERE statistic has been proposed as a similarity index between RNA-Seq samples [Schulze, 2012]. It measures whether the variability between samples is random Poisson variability or higher. Pairwise SERE values are printed in the lower triangle of the pairwise scatter plot. The value of the SERE statistic is:

- 0 when samples are identical (no variability at all: this may happen in the case of a sample duplication);
- 1 for technical replicates (technical variability follows a Poisson distribution);
- greater than 1 for biological replicates and samples from different biological conditions (biological variability is higher than technical one, data are over-dispersed with respect to Poisson). The higher the SERE value, the lower the similarity. It is expected to be lower between biological replicates than between samples of different biological conditions. Hence, the SERE statistic can be used to detect inversions between samples.

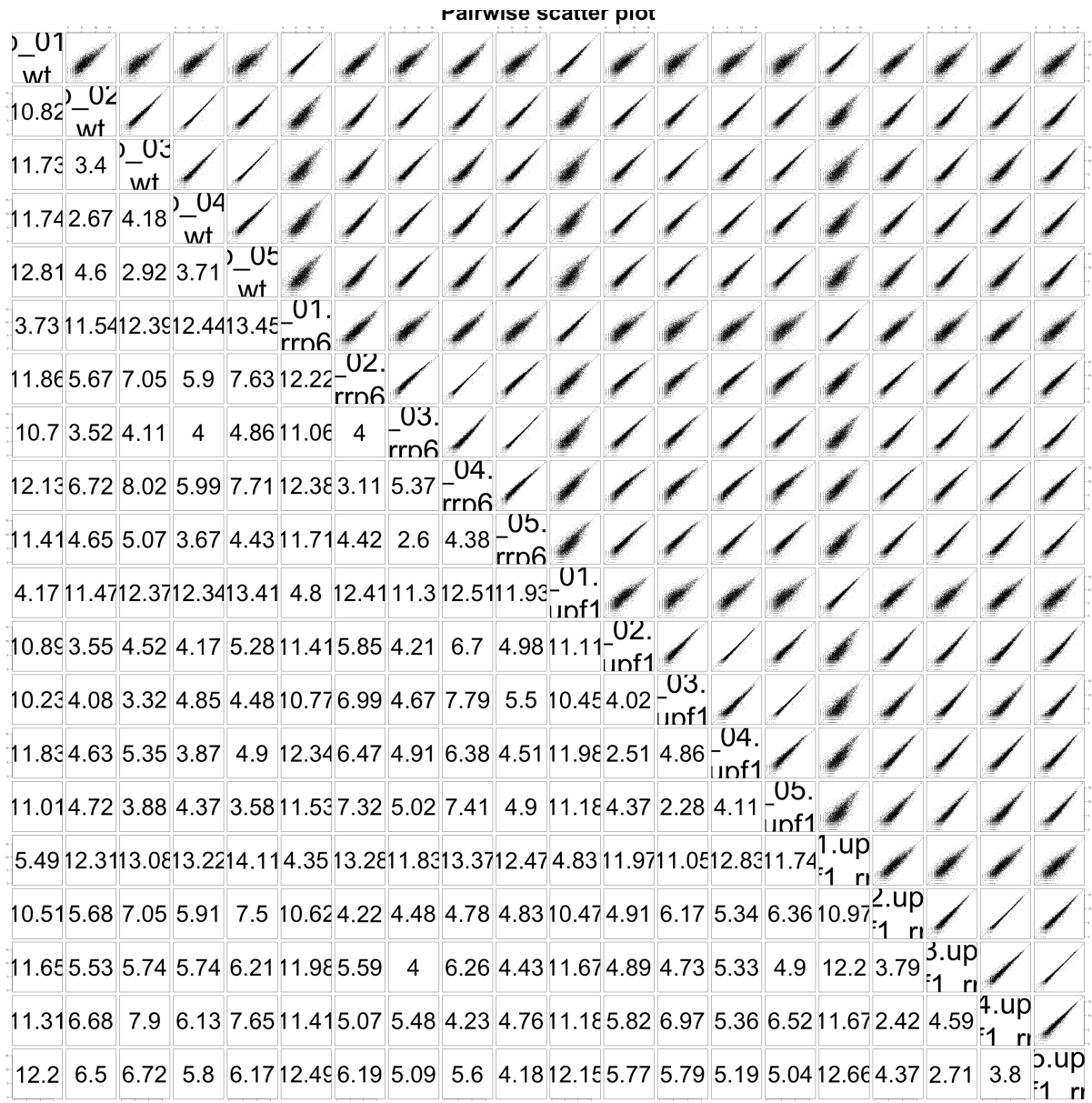


Figure 5: Pairwise comparison of samples.

3 Variability within the experiment: data exploration

The main variability within the experiment is expected to come from biological differences between the samples. This can be checked in two ways. The first one is to perform a hierarchical clustering of the whole sample set. This is performed after a transformation of the count data which can be either a Variance Stabilizing Transformation (VST) or a regularized log transformation (rlog) [Anders, 2010 and Love, 2014].

A VST is a transformation of the data that makes them homoscedastic, meaning that the variance is then independent of the mean. It is performed in two steps: (i) a mean-variance relationship is estimated from the data with the same function that is used to normalize count data and (ii) from this relationship, a transformation of the data is performed in order to get a dataset in which the variance is independent of the mean. The homoscedasticity is a prerequisite for the use of some data analysis methods, such as hierarchical clustering or Principal Component Analysis (PCA). The regularized log transformation is based on a GLM (Generalized Linear Model) on the counts and has the same goal as a VST but is more robust in the case when the size factors vary widely.

Figure 6 shows the dendrogram obtained from VST-transformed data. An euclidean distance is computed between samples, and the dendrogram is built upon the Ward criterion. We expect this dendrogram to group replicates and separate biological conditions.

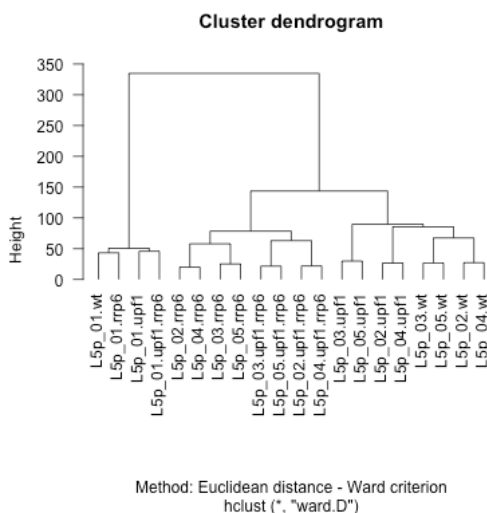


Figure 6: Sample clustering based on normalized data.

Another way of visualizing the experiment variability is to look at the first principal components of the PCA, as shown on the figure 7. On this figure, the first principal component (PC1) is expected to separate samples from the different biological conditions, meaning that the biological variability is the main source of variance in the data.

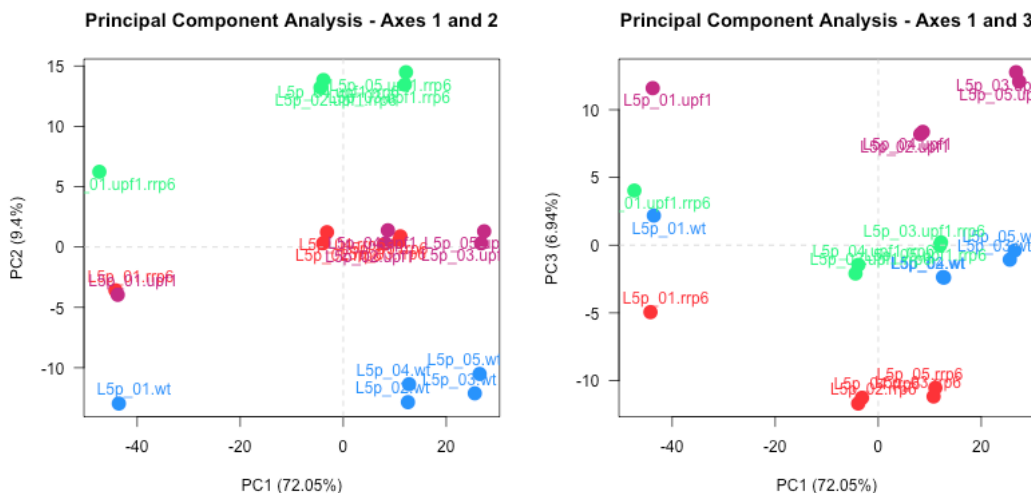


Figure 7: First two components of a Principal Component Analysis, with percentages of variance associated with each axis.

4 Normalization

Normalization aims at correcting systematic technical biases in the data, in order to make read counts comparable across samples. The normalization proposed by DESeq2 relies on the hypothesis that most features are not differentially expressed. It computes a scaling factor for each sample. Normalized read counts are obtained by dividing raw read counts by the scaling factor associated with the sample they belong to. Scaling factors around 1 mean (almost) no normalization is performed. Scaling factors lower than 1 will produce normalized counts higher than raw ones, and the other way around. Two options are available to compute scaling factors: `locfunc="median"` (default) or `locfunc="shorth"`. Here, the normalization was performed with `locfunc="shorth"`.

	LT_01.WT	LT_02.WT	LT_03.WT	LT_01.upf1	LT_02.upf1	LT_03.upf1	LT_01.upf1.xrn1	LT_02.upf1.xrn1	LT_03.upf1.xrn1	LT_01.xrn1
Size factor	1.10	0.62	1.73	0.96	0.66	1.62	1.01	0.61	1.61	0.9

Table 5: Normalization factors.

The histograms (figure 8) can help to validate the choice of the normalization parameter ("median" or "shorth"). Under the hypothesis that most features are not differentially expressed, each size factor is expected to be close to the mode of the distribution of the counts divided by their geometric means across samples.

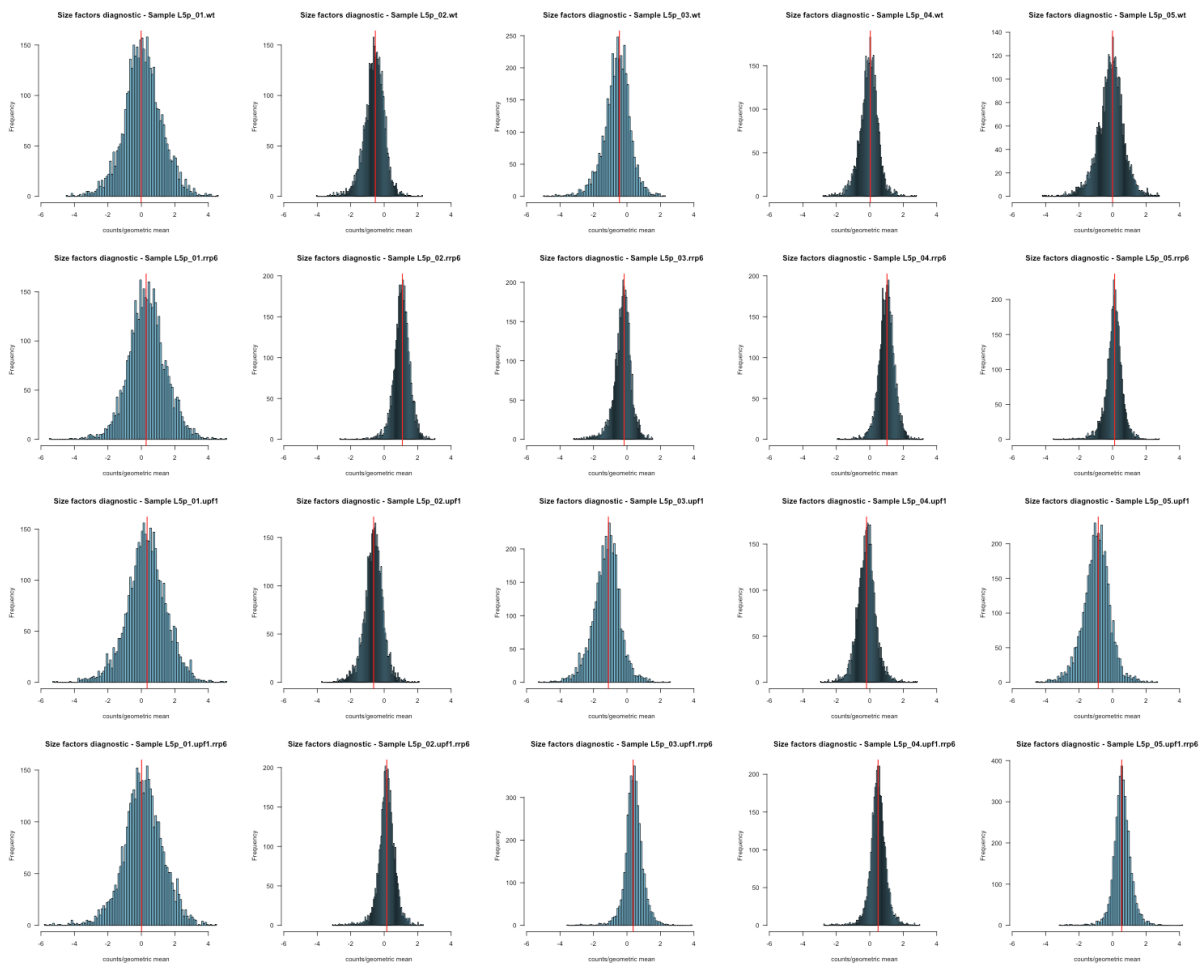


Figure 8: Diagnostic of the estimation of the size factors.

The figure 9 shows that the scaling factors of DESeq2 and the total count normalization factors may not perform similarly.

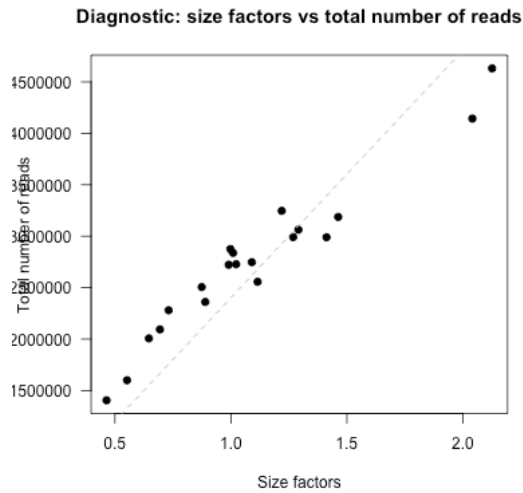


Figure 9: Plot of the estimated size factors and the total number of reads per sample.

Boxplots are often used as a qualitative measure of the quality of the normalization process, as they show how distributions are globally affected during this process. We expect normalization to stabilize distributions across samples. Figure 10 shows boxplots of raw (left) and normalized (right) data respectively.

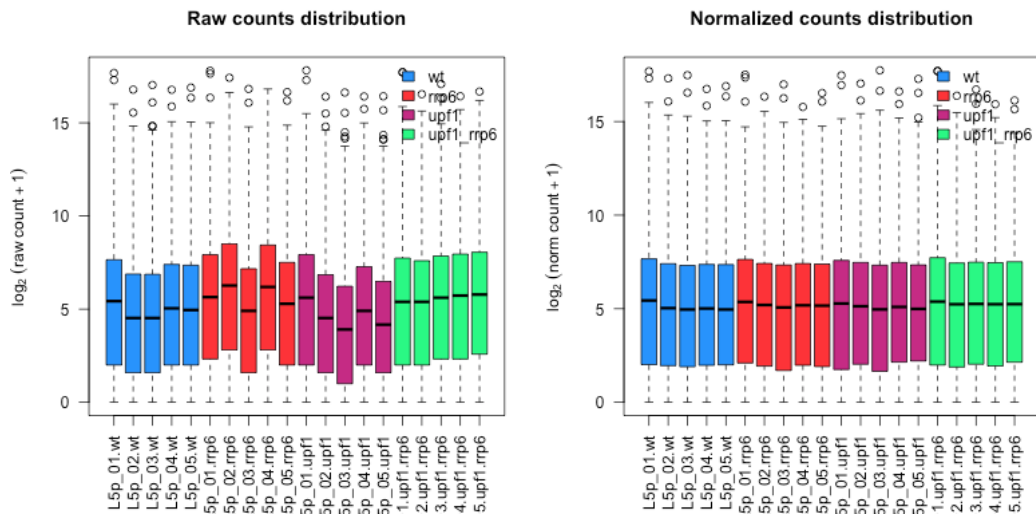


Figure 10: Boxplots of raw (left) and normalized (right) read counts.

5 Differential analysis

5.1 Modelisation

DESeq2 aims at fitting one linear model per feature. For this project, the design used is counts ~ group and the goal is to estimate the models' coefficients which can be interpreted as $\log_2(FC)$. These coefficients will then be tested to get p-values and adjusted p-values.

5.2 Outlier detection

Model outliers are features for which at least one sample seems unrelated to the experimental or study design. For every feature and for every sample, the Cook's distance [Cook, 1977] reflects how the sample matches the model. A large value of the Cook's distance indicates an outlier count and p-values are not computed for the corresponding feature.

5.3 Dispersions estimation

The DESeq2 model assumes that the count data follow a negative binomial distribution which is a robust alternative to the Poisson law when data are over-dispersed (the variance is higher than the mean). The first step of the statistical procedure is to estimate the dispersion of the data. Its purpose is to determine the shape of the mean-variance relationship. The default is to apply a GLM (Generalized Linear Model) based method (fitType="parametric"), which can handle complex designs but may not converge in some cases. The alternative is to use fitType="local" as described in the original paper [Anders, 2010]. The parameter used for this project is fitType="parametric". Then, DESeq2 imposes a Cox Reid-adjusted profile likelihood maximization [Cox, 1987 and McCarthy, 2012] and uses the maximum *a posteriori* (MAP) of the dispersion [Wu, 2013].

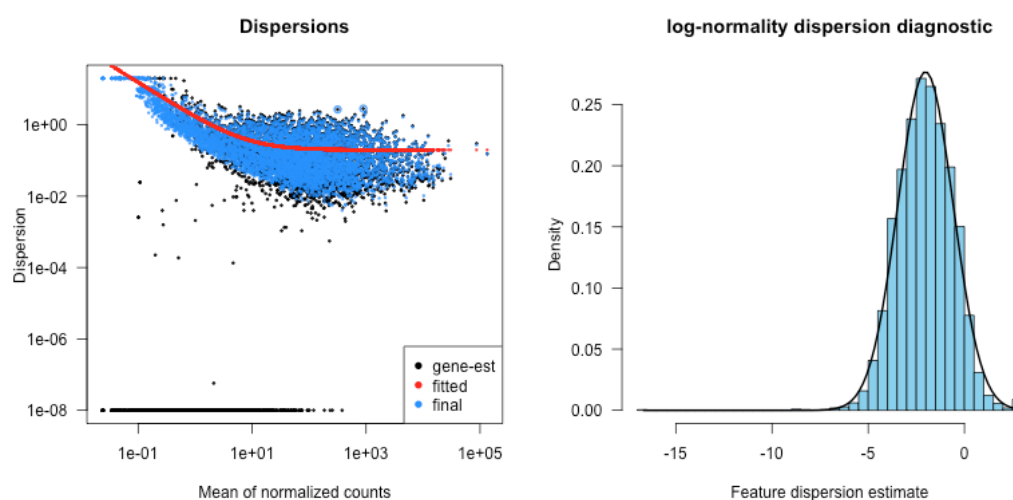


Figure 11: Dispersion estimates (left) and diagnostic of log-normality (right).

The left panel on figure 11 shows the result of the dispersion estimation step. The x- and y-axes represent the mean count value and the estimated dispersion respectively. Black dots represent empirical dispersion estimates for each feature (from the observed counts). The red dots show the mean-variance relationship function (fitted dispersion value) as estimated by the model. The blue dots are the final estimates from the maximum *a posteriori* and are used to perform the statistical test. Blue circles (if any) point out dispersion outliers. These are features with a very high empirical variance (computed from observed counts). These high dispersion values fall far from the model estimation. For these features, the statistical test is based on the empirical variance in order to be more conservative than with the MAP dispersion. These features will have low chance to be declared significant. The figure on the right panel allows to check the hypothesis of log-normality of the dispersions.

5.4 Statistical test for differential expression

Once the dispersion estimation and the model fitting have been done, DESeq2 can perform the statistical testing. Figure 12 shows the distributions of raw p-values computed by the statistical test for the comparison(s) done. This distribution is expected to be a mixture of a uniform distribution on $[0, 1]$ and a peak around 0 corresponding to the differentially expressed features.

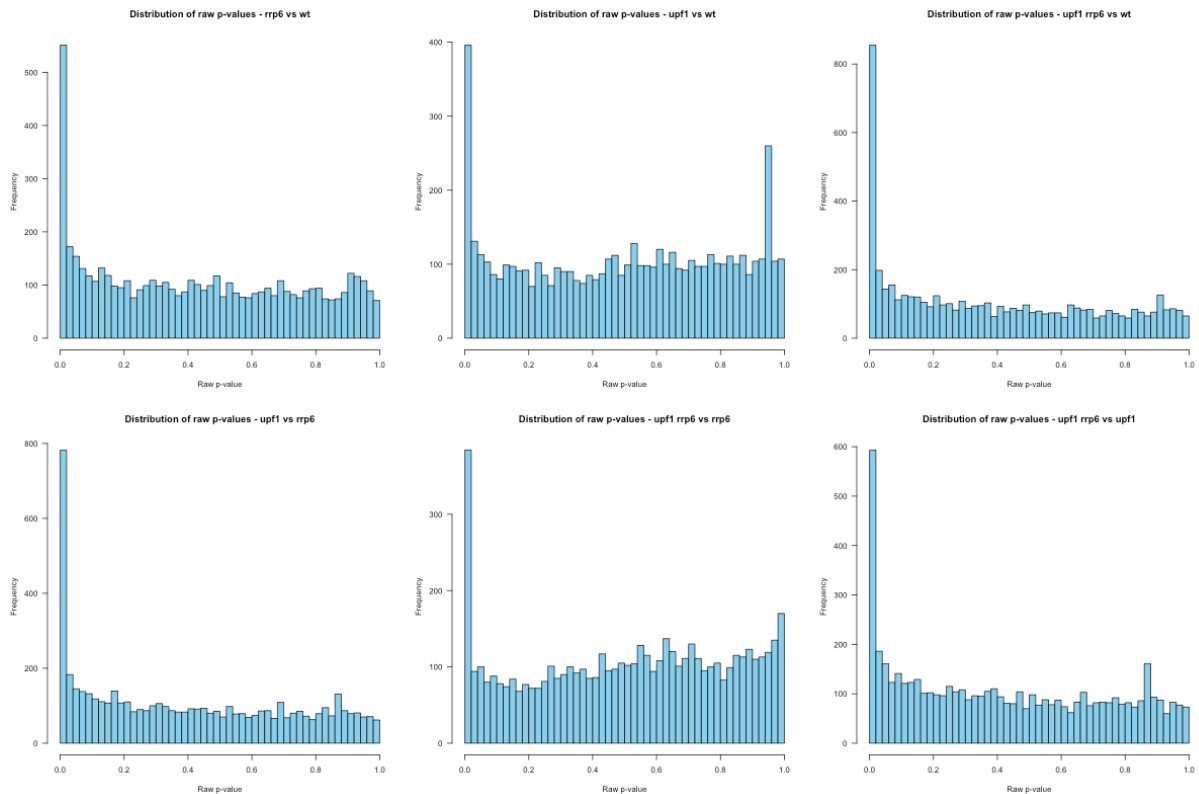


Figure 12: Distribution(s) of raw p-values.

5.5 Independent filtering

DESeq2 can perform an independent filtering to increase the detection power of differentially expressed features at the same experiment-wide type I error. Since features with very low counts are not likely to see significant differences typically due to high dispersion, it defines a threshold on the mean of the normalized counts irrespective of the biological condition. This procedure is independent because the information about the variables in the design formula is not used [Love, 2014].

Table 6 reports the thresholds used for each comparison and the number of features discarded by the independent filtering. Adjusted p-values of discarded features are then set to NA.

```
## Error in print(xtable(summaryResults$tabIndepFiltering, caption = "Table 6: Number of features discarded by
## pas de méthode pour 'xtable' applicable pour un objet de classe "NULL"
```

5.6 Final results

A p-value adjustment is performed to take into account multiple testing and control the false positive rate to a chosen level α . For this analysis, a BH p-value adjustment was performed [Benjamini, 1995 and 2001] and the level of controlled false positive rate was set to 0.05.

```
## Error in print(xtable(summaryResults$nDiffTotal, caption = paste0(ifelse(independentFiltering, : erreur d'év
## pas de méthode pour 'xtable' applicable pour un objet de classe "NULL"
```

Figure 13 represents the MA-plot of the data for the comparisons done, where differentially expressed features are highlighted in red. A MA-plot represents the log ratio of differential expression as a function of the mean intensity for each feature.

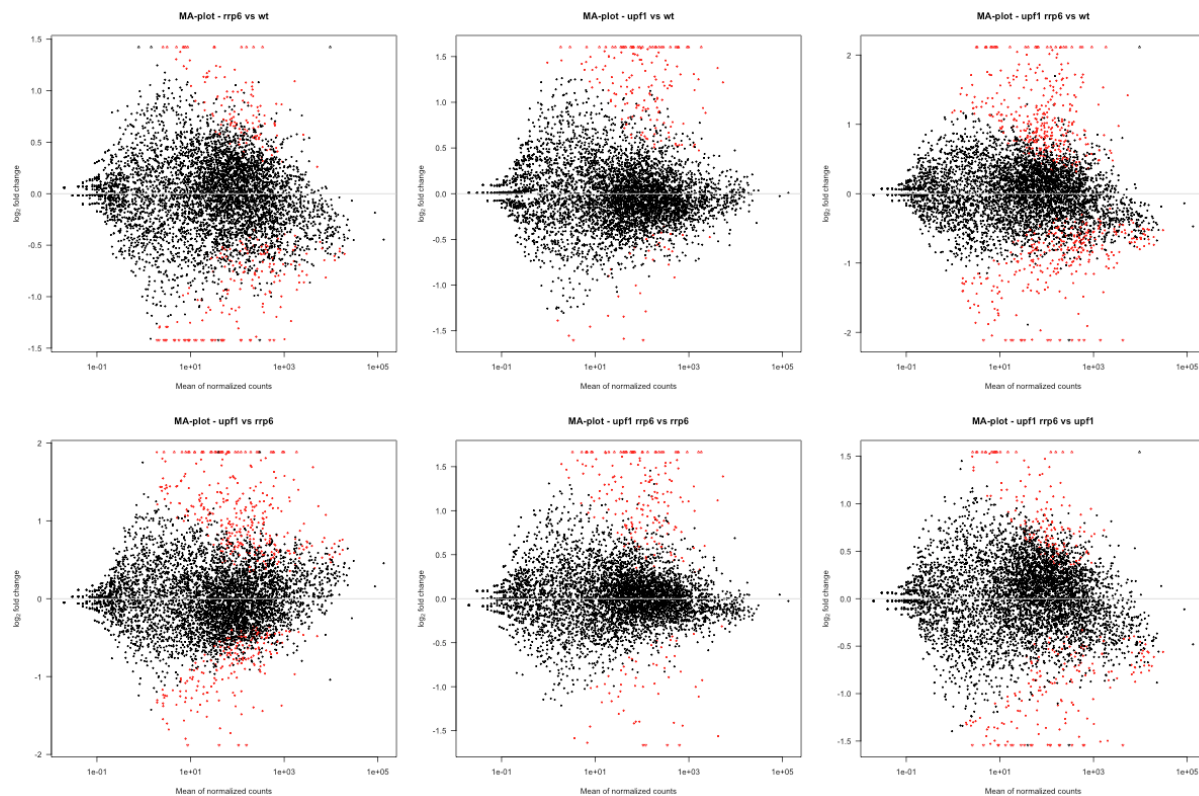


Figure 13: MA-plot(s) of each comparison. Red dots represent significantly differentially expressed features.

Figure 14 shows the volcano plots for the comparisons performed and differentially expressed features are still highlighted in red. A volcano plot represents the log of the adjusted P value as a function of the log ratio of differential expression.

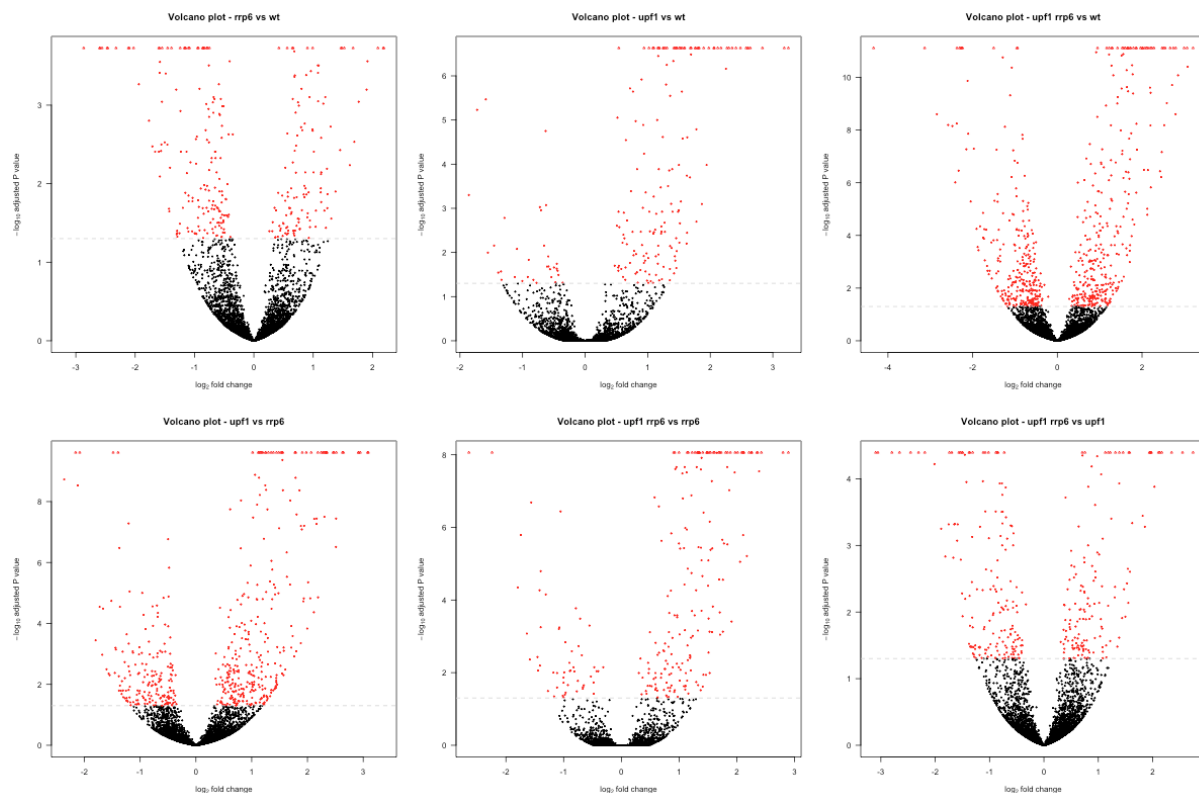


Figure 14: Volcano plot(s) of each comparison. Red dots represent significantly differentially expressed features.

Full results as well as lists of differentially expressed features are provided in the following text files which can be easily read in a spreadsheet. For each comparison:

- TestVsRef.complete.txt contains results for all the features;
- TestVsRef.up.txt contains results for significantly up-regulated features. Features are ordered from the most significant adjusted p-value to the less significant one;
- TestVsRef.down.txt contains results for significantly down-regulated features. Features are ordered from the most significant adjusted p-value to the less significant one.

These files contain the following columns:

- Id: unique feature identifier;
- sampleName: raw counts per sample;
- norm.sampleName: rounded normalized counts per sample;
- baseMean: base mean over all samples;
- WT, upf1, upf1_xrn1 and xrn1: means (rounded) of normalized counts of the biological conditions;
- FoldChange: fold change of expression, calculated as $2^{\log_2(\text{FC})}$;
- log2FoldChange: $\log_2(\text{FC})$ as estimated by the GLM model. It reflects the differential expression between Test and Ref and can be interpreted as $\log_2\left(\frac{\text{Test}}{\text{Ref}}\right)$. If this value is:
 - around 0: the feature expression is similar in both conditions;
 - positive: the feature is up-regulated (Test > Ref);
 - negative: the feature is down-regulated (Test < Ref);
- pvalue: raw p-value from the statistical test;
- padj: adjusted p-value on which the cut-off α is applied;
- dispGeneEst: dispersion parameter estimated from feature counts (i.e. black dots on figure 11);
- dispFit: dispersion parameter estimated from the model (i.e. red dots on figure 11);
- dispMAP: dispersion parameter estimated from the Maximum *A Posteriori* model;
- dispersion: final dispersion parameter used to perform the test (i.e. blue dots and circles on figure 11);
- betaConv: convergence of the coefficients of the model (TRUE or FALSE);
- maxCooks: maximum Cook's distance of the feature;
- outlier: indicates if the feature has been detected as a count outlier (it does not make sense if an infinite threshold is given for the Cook's distance).

6 R session information and parameters

The versions of the R software and Bioconductor packages used for this analysis are listed below. It is important to save them if one wants to re-perform the analysis in the same conditions.

R version 3.1.0 (2014-04-10) Platform: x86_64-apple-darwin13.1.0 (64-bit)

locale: [1] fr_FR.UTF-8/fr_FR.UTF-8/fr_FR.UTF-8/C/fr_FR.UTF-8/fr_FR.UTF-8

attached base packages: [1] parallel stats4 stats graphics grDevices utils datasets [8] grid methods base

other attached packages: [1] knitr_1.9 genefilter_1.48.1
 [3] SARTools_1.0.4 xtable_1.7-4
 [5] edgeR_3.8.5 limma_3.22.6
 [7] DESeq2_1.6.3 RcppArmadillo_0.4.650.1.1 [9] Rcpp_0.11.4 GenomicRanges_1.18.4
 [11] GenomeInfoDb_1.2.4 IRanges_2.0.1
 [13] S4Vectors_0.4.0 BiocGenerics_0.12.1
 [15] seqLogo_1.32.1

loaded via a namespace (and not attached): [1] acepack_1.3-3.3 annotate_1.44.0 AnnotationDbi_1.28.1 [4] base64enc_0.1-2

BatchJobs_1.5 BBmisc_1.9
 [7] Biobase_2.26.0 BiocParallel_1.0.3 brew_1.0-6
 [10] checkmate_1.5.1 cluster_2.0.1 codetools_0.2-10
 [13] colorspace_1.2-5 DBI_0.3.1 digest_0.6.8
 [16] evaluate_0.5.5 fail_1.2 foreach_1.4.2
 [19] foreign_0.8-63 formatR_1.0 Formula_1.2-0
 [22] geneplotter_1.44.0 ggplot2_1.0.0 gtable_0.1.2
 [25] Hmisc_3.15-0 iterators_1.0.7 lattice_0.20-30
 [28] latticeExtra_0.6-26 locfit_1.5-9.1 markdown_0.7.4
 [31] MASS_7.3-39 mime_0.2 munsell_0.4.2
 [34] nnet_7.3-9 plyr_1.8.1 proto_0.3-10
 [37] RColorBrewer_1.1-2 reshape2_1.4.1 rpart_4.1-9
 [40] RSQLite_1.0.0 scales_0.2.4 sendmailR_1.2-1
 [43] splines_3.1.0 stringr_0.6.2 survival_2.38-1
 [46] tcltk_3.1.0 tools_3.1.0 XML_3.98-1.1
 [49] XVector_0.6.0

Parameter values used for this analysis are:

- workDir: ~/bioinfo/bin/python/workspaces/nmd_effect/data/round2/data_for_publication/
- projectName: TSSC_nam7_xrn1
- author: Christophe Malabat
- targetFile: design_xrn1.txt
- rawDir: ~/bioinfo/bin/python/workspaces/nmd_effect/data/round2/data_for_publication/
- featuresToRemove: NULL
- varInt: group
- condRef: WT
- batch: NULL
- fitType: parametric
- cooksCutoff: TRUE
- independentFiltering: TRUE
- alpha: 0.05
- pAdjustMethod: BH
- typeTrans: VST
- locfunc: shorth
- colors: dodgerblue, firebrick1, MediumVioletRed, SpringGreen

7 Bibliography

- R Core Team, R: A Language and Environment for Statistical Computing, *R Foundation for Statistical Computing*, 2014
- Gentleman, Carey, Bates et al, Bioconductor: Open software development for computational biology and bioinformatics, *Genome Biology*, 2004
- Anders and Huber, Differential expression analysis for sequence count data, *Genome Biology*, 2010
- Love, Huber and Anders, Moderated estimation of fold change and dispersion for RNA-Seq data with DESeq2, *Genome Biology*, 2014
- Dillies, Rau, Aubert et al, A comprehensive evaluation of normalization methods for Illumina RNA-seq data analysis, *Briefings in Bioinformatics*, 2012
- Schulze, Kanwar, Golzenleuchter et al, SERE: Single-parameter quality control and sample comparison for RNA-Seq, *BMC Genomics*, 2012
- Cook, Detection of Influential Observation in Linear Regression, *Technometrics*, 1977
- Cox and Reid, Parameter orthogonality and approximate conditional inference, *Journal of the Royal Statistical Society*, 1987
- McCarthy, Chen and Smyth, Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation, *Nucleic Acids Research*, 2012
- Wu, Wang and Wu, A new shrinkage estimator for dispersion improves differential expression detection in RNA-seq data, *Biostatistics*, 2013
- Benjamini and Hochberg, Controlling the False Discovery Rate : A Practical and Powerful Approach to Multiple Testing, *Journal of the Royal Statistical Society*, 1995
- Benjamini and Yekutieli, The control of the false discovery rate in multiple testing under dependency, *The Annals of Statistics*, 2001