

Analysis of differential gene expression in wild and cultivated rice under drought stress

Ingo Giebel

Math.-Nat. Fakultät, Heinrich-Heine-Universität Düsseldorf

QBio304: Applied Bioinformatics

Prof. Dr. Björn Usadel

Dr. Jędrzej Jakub Szymanski

May 1, 2023

Contents

1 Abstract	3
2 Introduction	3
2.1 Background	3
2.2 Objectives	3
3 Materials and methods	3
3.1 Selection of the RNA-seq data	3
3.2 Quality evaluation	3
3.3 RNA-seq preprocessing	4
3.4 Transcripts abundances quantification	4
3.5 Statistical evaluation and differential expression analysis	5
3.5.1 Import of the kallisto transcript-level estimates	5
3.5.2 Filtering and normalization	6
3.5.3 Hierarchical cluster analysis	7
3.5.4 Principal component analysis	7
3.5.5 Identification of differentially expressed genes	9
3.6 Functional enrichment analysis	9
4 Results	9
4.1 Quality evaluation	9
4.2 Mapping efficiency and coverage	9
4.3 Exploratory data analysis	10
4.4 Differentially expressed genes	10
4.5 Functional enrichment analysis	10
5 Discussion	10
5.1 Critical evaluation of the results	10
5.2 Biological implications	10
5.3 Limitations and future directions	10
6 Conclusion	10
References	11
Index	13

List of Figures

1	TPM statistics about the imported kallisto data	6
2	Log2(CPM) distribution of the unfiltered, non-normalized data	6
3	Log2(CPM) distribution of the filtered (< 1 CPM in at least half of the samples), non-normalized data	7
4	Log2(CPM) distribution of the filtered, normalized data	7
5	Log2(CPM) distribution of the filtered, normalized data in comparison with the non-normalized and the unfiltered data	8

List of Tables

1	Number of genes with no reads at all, and conversely number of genes with CPMs ≥ 1 in at least $n = 1, 2, 3, \dots$ of the samples	6
---	---	---

1 Abstract

Provide a brief summary of the purpose of the assignment, the methods used, the main findings, and the significance of the results. Limit the abstract to 200-250 words.

2 Introduction

2.1 Background

Introduce and explain the study and give a rationale for the RNA-seq analysis. Discuss why of RNA sequencing in understanding gene expression and regulation in plants and why it is used in this study.

2.2 Objectives

State the specific aims of the assignment, which include obtaining plant RNA-seq data, evaluating its quality, mapping to a respective genome, performing statistical evaluation, differential expression analysis, functional enrichment analysis, and critically evaluating and discussing the results.

3 Materials and methods

3.1 Selection of the RNA-seq data

This study uses publicly available paired-end RNA-seq data of wild and cultivated rice, submitted in January 1, 2021 by the Institute of Botany, Chinese Academy of Sciences. This data allows to compare rice grown under normal conditions with rice grown under drought stress conditions. Furthermore, the data allows for an interspecies comparison of wild rice (*Oryza nivara*, cultivars BJ278 and BJ89) with cultivated rice (*Oryza sativa*, cultivar Nipponbare).

All samples were uniformly taken from seedlings (leaf tissue) at the age of twelve days. Used sequencing platform: Illumina HiSeq 2000.

Therefore, the data is well-suited for a targeted analysis of drought stress responses.

3.2 Quality evaluation

The quality of the raw and trimmed RNA-seq data was assessed using FastQC ("Babraham Bioinformatics," 2023). FastQC is a quality control analysis tool for high throughput sequencing data. It provides information about

- basic statistics: some simple composition statistics for the FastQ file analyzed

- per base sequence quality: an overview of the range of quality values across all bases at each position in the FastQ file
- per tile sequence quality: an overview of the per tile sequence quality in case an Illumina library was used
- per sequence quality scores: an overview of how the overall quality scores of the sequences are distributed
- per base sequence content: an overview of the proportion of each base position in a FastQ file for which each of the four normal DNA bases has been called
- per sequence GC content: the GC content across the whole length of each sequence in a file compared with a normal distributed GC content
- per base N content: an overview of the N content at each position across all bases
- sequence length distribution: an overview of how the sequence lengths are distributed
- sequence duplication levels: an overview of the degree of duplication for every sequence in a library
- over-represented sequences: a list of over-represented sequences matched against common contaminants
- adapter content: a check for significant amounts of adapter sequences the FastQ file

The results of the separate FastQC analyses (of the raw and trimmed FastQ files), the results of the Trimmomatic trimming and the information about the kallisto pseudoalignments were summarized in an interactive MultiQC HTML-report. See (Ewels et al., 2016).

3.3 RNA-seq preprocessing

The initial quality assessment of the raw FastQ files revealed that about roughly the first 12 base pairs of the reads were of low quality. Therefore, the data was preprocessed/trimmed using Trimmomatic, a "flexible and efficient preprocessing tool, which could correctly handle paired-end data" (Ewels et al., 2016).

3.4 Transcripts abundances quantification

The mapping/pseudoalignment of the RNA-seq reads and the abundances quantification of the transcripts was done using kallisto. According to (Bray et al., 2016), kallisto offers the following advantages over other alignment and quantification software:

kallisto is a program for quantifying abundances of transcripts from RNA-Seq data, or more generally of target sequences using high-throughput sequencing reads. It is based on the novel idea of pseudoalignment for rapidly determining the compatibility of reads with targets, without the need for alignment. On benchmarks with standard RNA-Seq data, kallisto can quantify 30 million human bulk RNA-seq reads in less than 3 minutes on a Mac desktop computer using only the read sequences and a transcriptome index that itself takes than 10 minutes to build. Pseudoalignment of reads preserves the key information needed for quantification, and kallisto is therefore not only fast, but also comparably accurate to other existing quantification tools. In fact, because the pseudoalignment procedure is robust to errors in the reads, in many benchmarks kallisto significantly outperforms existing tools.

kallisto requires a reference genome/transcriptome for aligning the RNA-seq data. To this end, reference FASTA cDNA dumps of *Oryza nivara*, cultivar BJ278 and *Oryza sativa*, cultivar Nipponbare were downloaded from Ensembl. The Ensembl project delivers reference data for genome interpretation for any species: genome assemblies from public archive are annotated with genes, regulatory regions, variants and comparative data to provide a foundation for scientific research and genome interpretation (Bolger et al., 2014).

3.5 Statistical evaluation and differential expression analysis

The preprocessed and aligned data was further evaluated and analyzed using an R-script (R Core Team, 2023) executed within RStudio (Posit team, 2023). The following sections describe the analysis steps in detail.

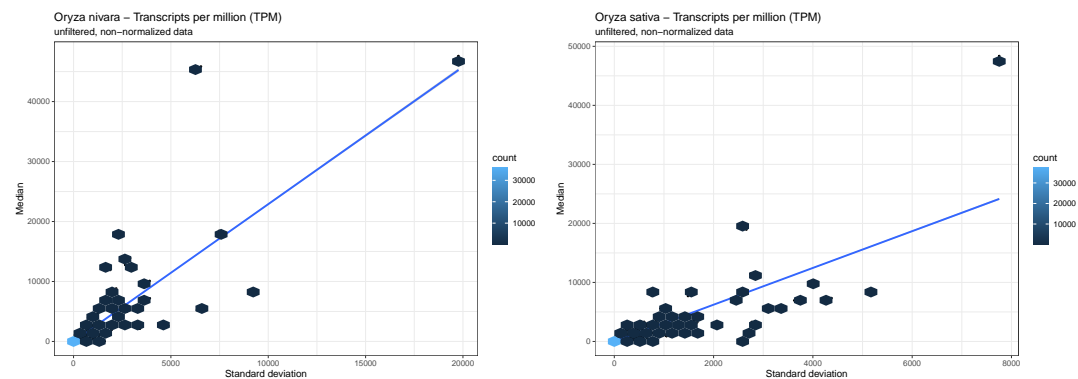
3.5.1 Import of the kallisto transcript-level estimates

The kallisto transcript-level estimates were imported using the R-package `tximport` (Love et al., 2022; Soneson et al., 2015). Thereby, the abundances, counts, and transcript lengths were summarized to the gene level.

Scaling method: average transcript length over samples and then the library size (parameter `lengthScaledTPM`).

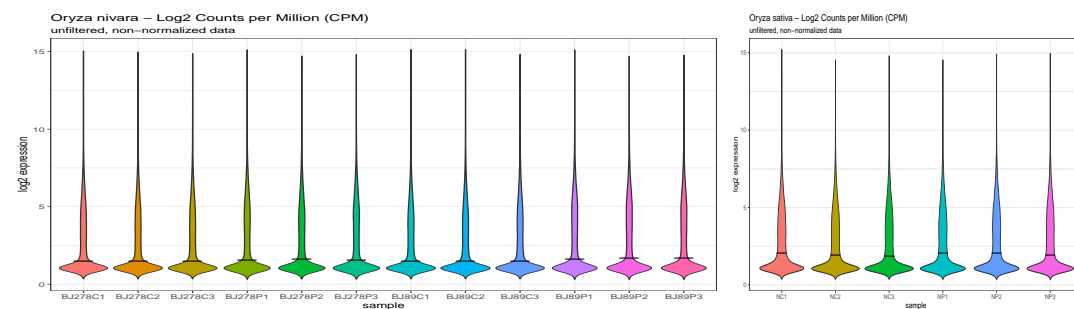
The mapping of the transcript IDs (used within the kallisto `abundance.tsv` files) to the corresponding gene IDs was done using the BioMart database `plants_mart` hosted at <https://plants.ensembl.org>. Datasets: `nivara_eg_gene` and `osativa_eg_gene` for *O. nivara* and *O. sativa*, respectively. See (Durinck & Huber, 2023; Durinck et al., 2009).

Figure 1 shows some basic transcripts per million (TPM) statistics about the imported kallisto files.

Figure 1: TPM statistics about the imported kallisto data

3.5.2 Filtering and normalization

For further analysis, `DGEList`-objects with counts per million (CPM) and $\log_2(\text{CPM})$ values were created using the R-package `edgeR` (Chen et al., 2023; Robinson et al., 2010). Figure 2 shows the distribution of the $\log_2(\text{CPM})$ values.

Figure 2: $\log_2(\text{CPM})$ distribution of the unfiltered, non-normalized data

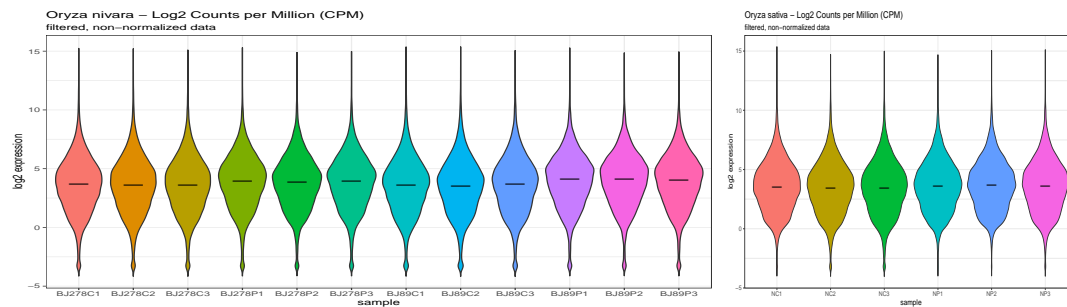
In order to assess reasonable values for filtering the data, the number of genes with no reads at all (in none of the samples) and conversely the number of genes with CPMs ≥ 1 in at least 1, 2, 3, ... of the samples were computed. Table 1 summarizes the results.

Table 1: Number of genes with no reads at all, and conversely number of genes with CPMs ≥ 1 in at least $n = 1, 2, 3, \dots$ of the samples

Species	Σ genes	Σ no reads	Genes with CPMs ≥ 1 in at least n samples					
			1	2	3	4	5	6
O. nivara	36313	7115	21001	20355	19891	19302	18895	18481
O. sativa	37967	4699	23934	22510	21609	20592	19715	18656

Filtering out genes with low reads (< 1 CPM in at least half of the samples) resulted in the distribution of the $\log_2(\text{CPM})$ values shown in figure 3.

Figure 3: $\log_2(\text{CPM})$ distribution of the filtered (< 1 CPM in at least half of the samples), non-normalized data



Finally, the filtered data was normalized using the `edgeR` function `calcNormFactors` which calculates scaling factors to convert raw library sizes into effective library sizes. Used normalization method: TMM. The results of the normalization are shown in figure 4.

Figure 4: $\log_2(\text{CPM})$ distribution of the filtered, normalized data

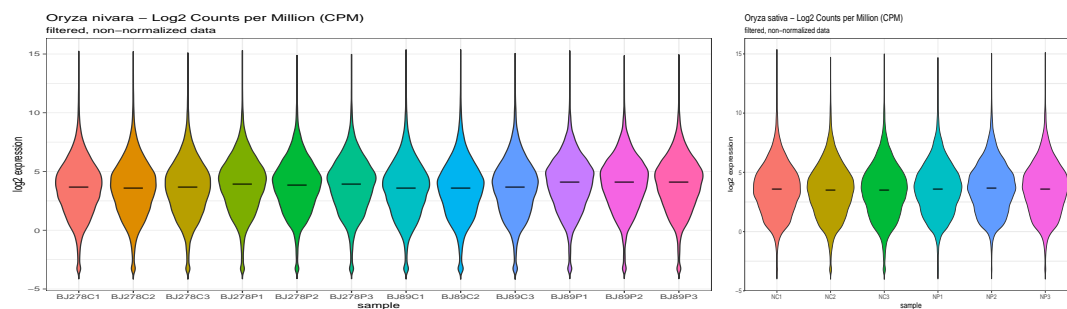


Figure 5 provides an overview of the filtering and normalization results.

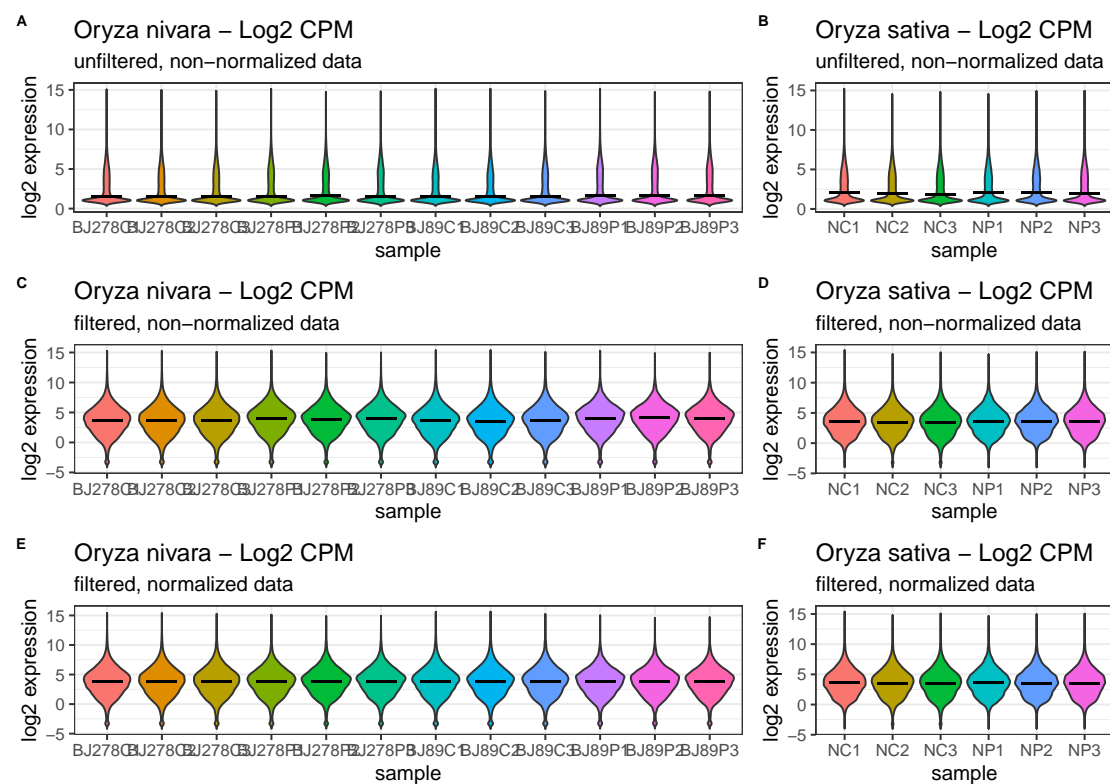
3.5.3 Hierarchical cluster analysis

A hierarchical cluster analysis (HCA) was performed via the `stats` function `hclust` on the Euclidean distance matrix of the $\log_2(\text{CPM})$ values. Used agglomeration method: `method = complete`.

3.5.4 Principal component analysis

A principal component analysis (PCA) was performed via the `stats` function `prcomp`.

Figure 5: Log2(CPM) distribution of the filtered, normalized data in comparison with the non-normalized and the unfiltered data



3.5.5 Identification of differentially expressed genes

In order to identify differentially expressed genes (DEGs), design and contrast matrices were created using the `stats` function `model.matrix` and the `limma` function `makeContrasts`.

The design matrices were used to create linear model fits for each gene via the `limma` functions `voom` and `lmFit`.

The linear model fits and the contrast matrices were used to calculate estimated coefficients and standard errors (contrasts) via the `limma` function `contrasts.fit`.

The contrasts were used to "calculate moderated t-statistics, moderated F-statistic, and log-odds of differential expression by empirical Bayes moderation of the standard errors towards a global value" (empirical Bayes statistics) via the `limma` function `eBayes` (Smyth et al., 2023).

Finally, the empirical Bayes statistics were used to extract a table of the top-ranked genes via the `limma` function `topTable` (sorted by the LogFC-values).

Venn diagrams of the DEGs (according to the calculated empirical Bayes statistics) were created via the `gprofiler2` function `decideTests` (with parameters `method = "global"`, `adjust.method = "BH"`, `p.value = 0.01`, `lfc = 7`) and the `limma` function `vennDiagram`.

See (Ritchie et al., 2015) for details on the statistical foundations implemented by `limma`.

3.6 Functional enrichment analysis

A functional enrichment analysis of the 100 "top-ranked" genes was performed via the `gprofiler2` function `gost` (with `correction_method = "fdr"` (Kolberg & Raudvere, 2021)).

4 Results

4.1 Quality evaluation

Present the findings from the quality evaluation of the selected RNA-seq data.

4.2 Mapping efficiency and coverage

Report the results of the mapping process, including the mapping efficiency and coverage.

4.3 Exploratory data analysis

Discuss the dominating variance components, reproducibility, possible batch effects and confounding variables.

4.4 Differentially expressed genes

Discuss the identified differentially expressed genes and their potential biological significance.

4.5 Functional enrichment analysis

Present the results of the functional enrichment analysis, highlighting the enriched functional categories.

5 Discussion

5.1 Critical evaluation of the results

Discuss the quality and reliability of the RNA-seq data and the downstream analyses.

5.2 Biological implications

Discuss the potential implications of the findings for plant biology and the broader scientific community.

5.3 Limitations and future directions

Address the limitations of the current analysis and suggest possible future directions to expand on the findings.

6 Conclusion

Summarize the main findings of the assignment, reiterating the significance of the results, and provide a final statement on the overall outcome of the study.

References

- Babraham bioinformatics. (2023, April 30). <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
- Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*, 30(15), 2114–2120. <https://doi.org/10.1093/bioinformatics/btu170>
- Bray, N. L., Pimentel, H., Melsted, P., & Pachter, L. (2016). Near-optimal probabilistic RNA-seq quantification. *Nature Biotechnology*, 34(5), 525–527. <https://doi.org/10.1038/nbt.3519>
- Chen, Y., Lun, A. T., McCarthy, D. J., Ritchie, M. E., Phipson, B., Hu, Y., Zhou, X., Robinson, M. D., & Smyth, G. K. (2023). *Edger: Empirical analysis of digital gene expression data in r* [<http://bioinf.wehi.edu.au/edgeR>].
- Durinck, S., & Huber, W. (2023). *Biomart: Interface to biomart databases (i.e. ensembl)* [R package version 2.54.1].
- Durinck, S., Spellman, P. T., Birney, E., & Huber, W. (2009). Mapping identifiers for the integration of genomic datasets with the r/bioconductor package biomart. *Nature Protocols*, 4, 1184–1191.
- Ewels, P., Magnusson, M., Lundin, S., & Käller, M. (2016). MultiQC: summarize analysis results for multiple tools and samples in a single report. *Bioinformatics*, 32(19), 3047–3048. <https://doi.org/10.1093/bioinformatics/btw354>
- Kolberg, L., & Raudvere, U. (2021). *Gprofiler2: Interface to the g:profiler toolset* [R package version 0.2.1]. <https://CRAN.R-project.org/package=gprofiler2>
- Love, M., Soneson, C., Robinson, M., Patro, R., Morgan, A. P., Thompson, R. C., Shirley, M., & Srivastava, A. (2022). *Tximport: Import and summarize transcript-level estimates for transcript- and gene-level analysis* [R package version 1.26.1]. <https://github.com/mikelove/tximport>
- Posit team. (2023). *Rstudio: Integrated development environment for r*. Posit Software, PBC. Boston, MA. <http://www.posit.co/>
- R Core Team. (2023). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. Vienna, Austria. <https://www.R-project.org/>
- Ritchie, M. E., Phipson, B., Wu, D., Hu, Y., Law, C. W., Shi, W., & Smyth, G. K. (2015). limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Research*, 43(7), e47. <https://doi.org/10.1093/nar/gkv007>
- Robinson, M. D., McCarthy, D. J., & Smyth, G. K. (2010). Edger: A bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*, 26(1), 139–140. <https://doi.org/10.1093/bioinformatics/btp616>
- Smyth, G., Hu, Y., Ritchie, M., Silver, J., Wettenhall, J., McCarthy, D., Wu, D., Shi, W., Phipson, B., Lun, A., Thorne, N., Oshlack, A., de Graaf, C., Chen, Y., Langaas, M., Ferkingstad, E., Davy, M., Pepin, F., & Choi, D. (2023). *Limma: Linear models for microarray data* [R package version 3.54.2]. <http://bioinf.wehi.edu.au/limma>

Soneson, C., Love, M. I., & Robinson, M. D. (2015). Differential analyses for rna-seq: Transcript-level estimates improve gene-level inferences. *F1000Research*, 4. <https://doi.org/10.12688/f1000research.7563.1>

Index

C

counts per million (CPM), [6](#)

cultivar

- BJ278, [3](#)
- BJ89, [3](#)
- Nipponbare, [3](#)

D

DNA

- cDNA, [5](#)

differentially expressed genes (DEGs), [9](#)

E

Ensembl, [5](#)

F

FASTA, [5](#)

FastQ, [3](#)

FastQC, [3](#)

G

gene

- ID, [5](#)

genome/transcriptome

- reference, [5](#)

H

hierarchical cluster analysis (HCA), [7](#)

- Euclidean distance matrix, [7](#)

I

Illumina

- HiSeq 2000, [3](#)

K

kallisto, [4](#)

L

LogFC, [9](#)

M

MultiQC, [4](#)

O

Oryza

- nivara, [3](#)
- sativa, [3](#)

P

principal component analysis (PCA), [7](#)

R

R, [5](#)

- RStudio, [5](#)
- package
 - edgeR, [6](#)
 - limma, [9](#)
 - stats, [7](#), [9](#)
 - tximport, [5](#)

RNA-seq, [3](#)

- adapter sequence, [4](#)
- contaminant, [4](#)
- quality assessment, [3](#)

T

transcript

- ID, [5](#)
- abundances quantification, [4](#)
- pseudoalignment, [4](#)
- transcripts per million (TPM), [5](#)

Trimmomatic, [4](#)