

Détection de gènes différentiellement exprimés à partir de données RNA-seq

Diplôme Interuniversitaire en Bioinformatique intégrative (DU-Bii 2019)

Jacques van Helden

2019-02-24

Contents

Introduction	1
But de ce TP	2
Cas d'étude	2
Sources des données	2
Paramètres de l'analyse	2
Téléchargement des données	3
Exploration des données	3
Statistiques descriptives	4
Distribution des comptages	4
Boîtes à moustache	5

Introduction

Depuis l'avènement du séquençage massivement parallèle (NGS, Next Generation Sequencing) en 2007, la détection de gènes différentiellement exprimés (DEG) à partir de données transcriptomiques RNA-seq constitue l'une de ses applications les plus populaires.

Le principe est de mesurer, dans différentes conditions, la concentration d'ARN correspondant à chaque gène, et de comparer ces concentrations entre les échantillons de deux (**comparaison à 2 groupes** ou **binaire**) ou plusieurs conditions (comparaison **multi-groupes**).

Dès les premières analyses, les chercheurs se sont rendus compte que les méthodes classiques d'analyse différentielles (tests de Student, ANOVA) ne se prêtaient pas du tout à l'analyse de telles données, car leur nature diffère fondamentalement des hypothèses de travail sous-jacentes, pour différentes raisons.

1. Le niveau d'expression d'un gène est quantifié sur une échelle discrète (comptage du nombre de lectures alignées qui chevauchent le gène), alors que les tests paramétriques reposent sur une hypothèse de normalité.
2. Les ordres de grandeur des comptages varient fortement d'un gène à l'autre: certains gènes ont une poignée de comptages par échantillon, d'autres des centaines de milliers. Les gènes très très fortement représentés correspondent généralement à des gènes non-codants (par exemple ARN ribosomique) ou à des fonctions cellulaires particulières liées au métabolisme de l'ARN, et on comprend les raisons biologiques de leur sur-représentation. Ils n'en constituent pas moins ce qu'on appelle en statistique des "valeurs aberrantes" ("outliers").
3. Les distributions de comptages comptent généralement un très grand nombre de valeurs nulles ("zero-inflated distributions").
4. Ces particularités posent des problèmes particuliers pour la normalisation des librairies de comptages :
 - la présence d'outliers affecte fortement la moyenne, et de façon très instable, ce qui biaise fortement l'estimation de la tendance centrale.

- la stratégie de repli sur des estimateurs robustes, comme la médiane est contestable du fait du très grand nombre de valeurs nulles (dans certains cas, plus d'un quart voire la moitié des gènes ont une valeur nulle pour un échantillon).

Des méthodes spécifiques ont donc été développées dès 2010 pour affronter ces difficultés particulières de la normalisation et de l'analyse différentielle des données de RNA-seq.

But de ce TP

Le but de ce TP est d'effectuer une première exploration de l'analyse différentielle des données d'expression, sur base d'un petit cas d'étude simple: l'analyse transcriptionnelle de mutants de sporulation chez la levure *Saccharomyces cerevisiae*.

Cas d'étude

... A ECRIRE: DESCRIPTION DU CAS D'ETUDE

Sources des données

Paramètres de l'analyse

Nous définissons dans une variable R (de type liste) les paramètres de l'analyse. Ceci nous permettra de reproduire ultérieurement exactement la même succession d'étapes soit en utilisant des données différentes, soit en modifiant les paramètres particuliers (design, seuils, ...).

```
## Load libraries
message("Loading libraries")
library(knitr)
library(kableExtra)
library(FactoMineR)
library(clues)
library(RColorBrewer)
library(ComplexHeatmap)

# library(corrplot)
library(FactoMineR)
library(ClassDiscovery)
# library(formattable)
```

```
parameters <- list(
  data.folder = "data/GSE89530", # dossier des données
  # data.folder = "/shared/projects/du_bii_2019/data/module3/seance5/GSE89530" # on the IFF-cluster-core
  counts = "GSE89530_counts.tsv.gz",
  sample.descr = "GSE89530_samples.tsv.gz",
  alpha = 0.05, # seuil de significativité
  epsilon = 0.1 # pseudo-comptage pour la transformation log2
)
```

```
kable(t(data.frame(parameters)), col.names = "Parameter Value")
```

	Parameter Value
data.folder	data/GSE89530
counts	GSE89530_counts.tsv.gz
sample.descr	GSE89530_samples.tsv.gz
alpha	0.05
epsilon	0.1

Table 1: ****Table de comptage de lectures (short reads) par gène.**** Sélection arbitraire de quelques gènes (lignes) et échantillons (colonnes)

	Condition	Replicate	Genotype	Time.point	Label
GSM2375291	Wild_type_4	3	Wild-type	4	Sc_WT_4h
GSM2375282	bdf1_Y187F_Y354F_mutant_4	1	bdf1-Y187F-Y354F mutant	4	bd1bd2mut_4
GSM2375281	bdf1_Y187F_Y354F_mutant_0	1	bdf1-Y187F-Y354F mutant	0	bd1bd2mut_0
GSM2375288	bdf1_Y187F_Y354F_mutant_4	2	bdf1-Y187F-Y354F mutant	4	bd1bd2mut_4

Téléchargement des données

```
## List files in the data folder
# list.files(parameters$data.folder)

## Load counts of reads per gene
counts <- read.delim(file.path(parameters$data.folder, parameters$counts), row.names = 1, header = TRUE)

## Load sample descriptions
samples <- read.delim(file.path(parameters$data.folder, parameters$sample.descr), row.names = 1, header = TRUE)
```

La table de comptages comporte 7478 lignes (correspondant aux gènes) et 18 colonnes (correspondant aux échantillons). Nous pouvons afficher un petit morceau de cette table, en sélectionnant au hasard quelques lignes et colonnes.

```
some.genes <- sample(1:nrow(counts), size = 10, replace = FALSE)
some.samples <- sample(1:ncol(counts), size = 4, replace = FALSE)
kable(counts[some.genes, some.samples])
```

	GSM2375289	GSM2375284	GSM2375278	GSM2375285
YER061C	390	355	653	337
YGL194C-A	29	11	41	33
YMR197C	676	235	489	647
YLR244C	611	377	835	440
YJL200C	1206	264	646	569
YMR194W_intron	18	9	27	5
YLL033W	115	26	37	109
YNL269W	102	7	5	121
YJR061W	256	74	191	195
YGL210W	690	141	362	762

Le tableau GSE89530_samples.tsv.gz fournit une description de chaque échantillon.

```
# print(samples[some.samples, ])
kable(samples[some.samples, ],
      caption = "**Table de comptage de lectures (short reads) par gène.** Sélection arbitraire de quelques gènes et échantillons",
      align = "c")
```

Exploration des données

Avant toute autre chose, il convient de mener une exploration préliminaire des données, afin de se familiariser avec leur distribution.

Statistiques descriptives

La fonction R `summary()` calcule des statistiques de base pour chaque colonne d'une matrice ou data frame. Nous imprimons ici un sous-ensemble de ces échantillons.

```
summary(counts[, some.samples])
```

GSM2375289	GSM2375284	GSM2375278	GSM2375285
Min. : 0.0	Min. : 0	Min. : 0	Min. : 0.0
1st Qu.: 34.0	1st Qu.: 15	1st Qu.: 32	1st Qu.: 30.2
Median : 248.0	Median : 135	Median : 297	Median : 254.0
Mean : 955.7	Mean : 521	Mean : 1127	Mean : 809.9
3rd Qu.: 626.0	3rd Qu.: 362	3rd Qu.: 775	3rd Qu.: 595.0
Max. : 578443.0	Max. : 340576	Max. : 722346	Max. : 472819.0

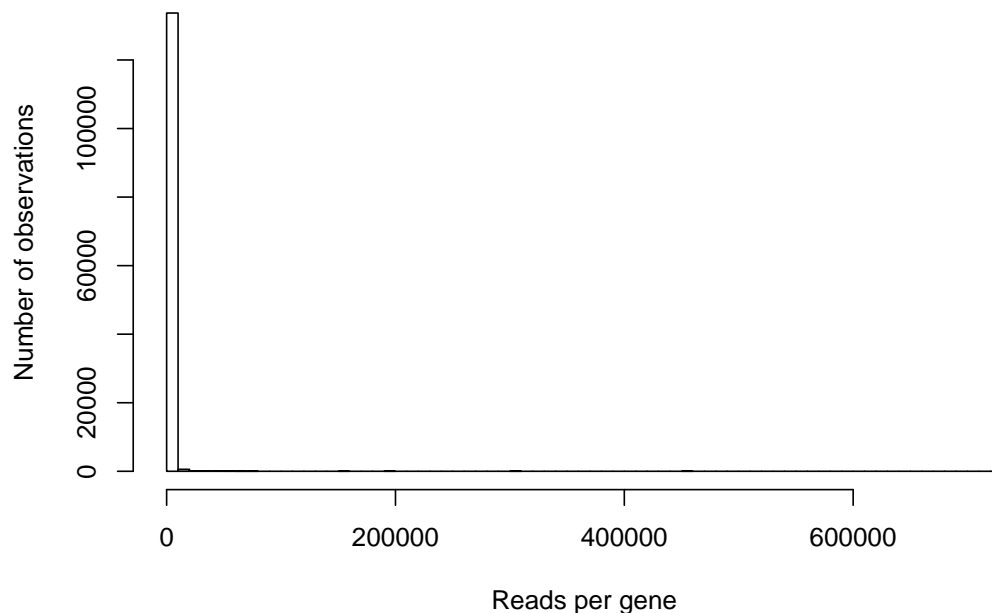
Distribution des comptages

```
count.stats <- mean(unlist(counts))
```

```
## Display histogram of the raw counts
```

```
hist(unlist(counts), breaks = 100, xlab = "Reads per gene", ylab = "Number of observations", main = "Distribution of raw counts")
```

Distribution of raw counts



Cet histogramme n'est pas très informatif, car toutes les valeurs sont concentrées dans la première tranche (à l'extrême gauche). Ceci résulte du fait que les intervalles de classe ont été définis sur base de l'étendue totale, et qu'il existe apparemment une observation qui a une valeur énorme par rapport aux autres (valeur aberrante, **"outlier"**). De fait, la valeur la plus élevée (722346) dépasse de très loin la moyenne (NA)

Afin de visualiser toute l'étendue des observations tout en mettant plus de détail sur les valeurs faibles, nous pouvons effectuer une transformation logarithmique. Nous devons cependant prêter attention à certains détails.

- Conventionnellement, on utilise la transformation \log_2 pour les données RNA-seq, car elle fournit un découpage plus fin des valeurs couvertes.

- Les données RnA-seq contiennent généralement un bon nombre de valeurs nulles (gènes non détectés), qui posent un problème pour la conversion logarithmique ($\log(0) = -\infty$). On contourne ce problème en ajoutant aux comptages un petit nombre arbitraire, qu'on dénomme **pseudo-comptage** et qu'on symbolise par la lettre grècque epsilon (ϵ).

Il est courant d'ajouter la valeur 1, mais nous préférons ajouter une valeur inférieure à 1, pour bien distinguer les comptages réels (qui peuvent de fait prendre une valeur 1) des pseudo-comptages.

Pour cette analyse nous choisissons $\epsilon = 0.1$.

La conversion donne donc.

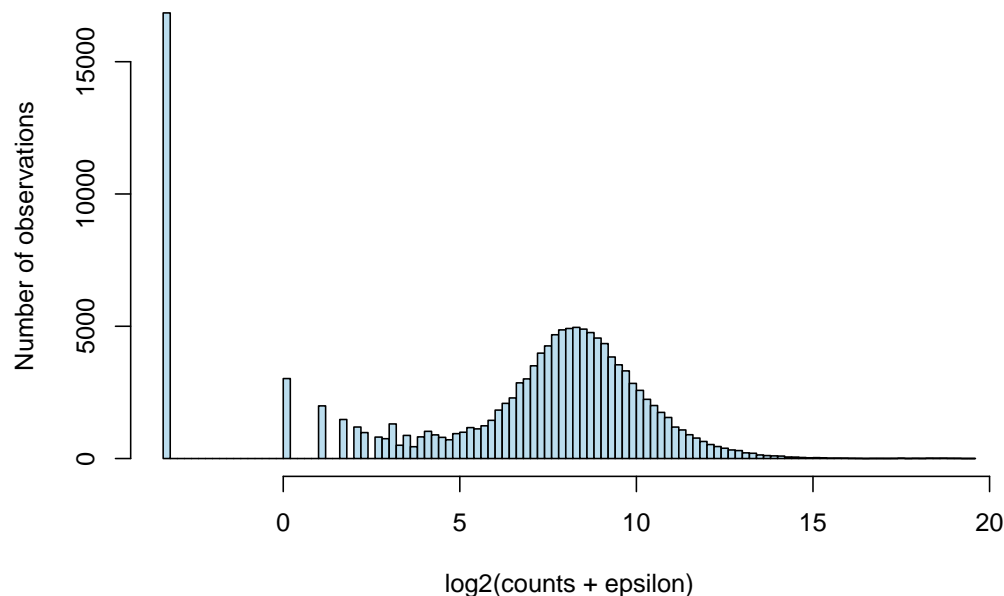
$$\log_2 \text{count} = \log_2(n + \epsilon) = \log_2(n + 0.1)$$

L'histogramme des $\log_2 \text{counts}$ est plus informatif que celui des comptages bruts.

```
counts.log2 <- log2(counts + parameters$epsilon)

hist(unlist(counts.log2),
     xlab = "log2(counts + epsilon)",
     ylab = "Number of observations",
     main = "Distribution of log2(counts)",
     breaks = 100, col = "#BBDDEE")
```

Distribution of $\log_2(\text{counts})$



Notons d'emblée le pic très élevé à gauche, qui correspond à toutes les observations nulles. Sur l'axe X , il apparaît à la valeur -3.3 , qui correspond à $\log_2(\epsilon) = \log_2(0.1)$.

Boîtes à moustache

Les boîtes à moustache sont très utilisées pour obtenir une vision d'ensemble de plusieurs distributions.

```
## Associate a specific color to each condition
conditions <- unique(samples$Condition)
condition.colors <- rainbow(n = length(conditions))
names(condition.colors) <- conditions
```

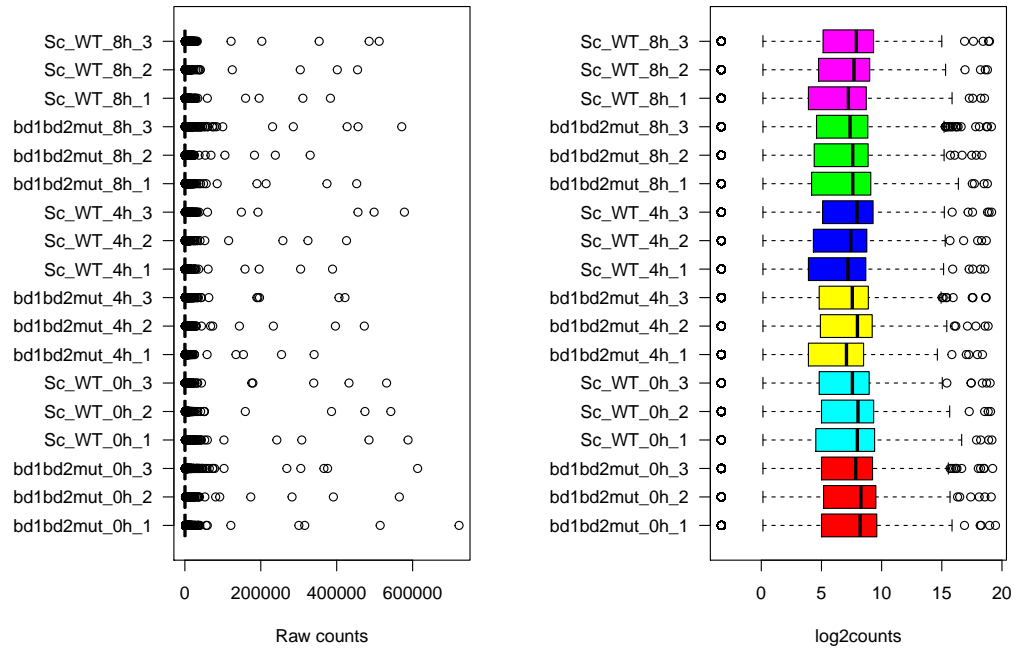


Figure 1: Boxplot of the counts per gene in each sample.

```
## Associate a color to each sample according to its condition
samples$color <- condition.colors[samples$Condition]

# Box plot of counts and log2-transformed counts
par(mfrow = c(1,2))
mar.ori <- par("mar") # Store the original value of the margin parameter
par(mar = c(5.1, 10.1, 4.1, 1.1)) # Increase leftmargin for sample labels

boxplot(counts, horizontal = TRUE, las = 1, names = samples$Label, col = samples$color, xlab = "Raw counts")
boxplot(counts.log2, horizontal = TRUE, las = 1, names = samples$Label, col = samples$color, xlab = "log2 counts")

par(mar = mar.ori) # reset original margins
par(mfrow = c(1,1))
```