



Statistical analysis of RNA-Seq data

G. Marot (Univ. Lille)

Sources: J. Aubert and C. Hennequet-Antier (Inrae)

M.A. Dillies and H. Varet (Institut Pasteur Paris)

Assistant for practical exercises: Samuel Blanck (Univ. Lille)

23 - 24 septembre 2021

Introduction

Differential analysis

Comparison of treatments, states, conditions, ...

Example : ill vs healthy

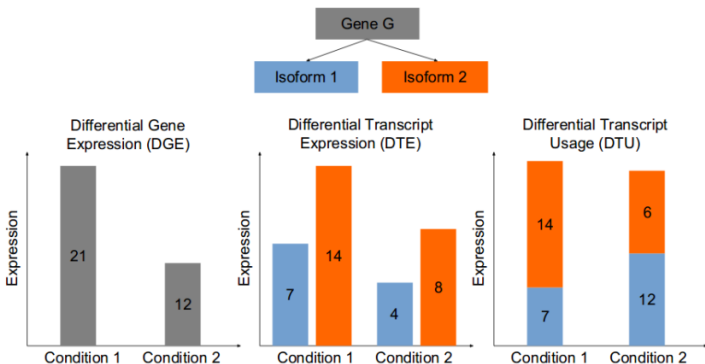
⇒ statistical analysis based on tests

Particularities of NGS data :

- Very few individuals
- Many tests (one per variable)
- Count data (statistical distributions different from the ones used for continuous data from microarrays)

Introduction

DGE : differential gene expression, DTE : differential transcript expression, DTU : differential transcript usage



This course focuses on DGE

Statistical test

- State the null and the alternative hypotheses
 - $H_0 = \{ \text{the mean expression (or proportion) of the gene is identical between the two conditions} \}$
 - $H_1 = \{ \text{the mean expression ((or proportion) of the gene is different between the two conditions} \}$
- Consider the statistical assumptions (e.g. independence) and distributions (e.g. normal, negative binomial, ...)
- Calculate the appropriate test statistic T
- Derive the distribution of the test statistic *under the null hypothesis* from the assumptions.
- Select a significance level (α), a probability threshold below which the null hypothesis will be rejected.

Remark : H_0 is always preferred. No sufficient proof \rightarrow no rejection.
 When we can not reject H_0 , this does not mean that H_0 is true.

Differential analysis

A gene is declared differentially expressed if the observed difference between two conditions is statistically significant, that is to say higher than some natural random variation.

Key steps for statisticians :

- experimental design
- normalization
- differential analysis
- multiple testing

Make an experimental design

Context of a RNA-seq experiment

Rule 0 : Share a common language in biology, bioinformatics and statistics.

Experimental design

All skills are needed to discussions right from project construction.

- **Rule 1** : Well define the biological question, get together and collect a priori knowledge (e.g. reference genome, splicing),
- **Rule 2** : Anticipate, identify all factors of variation and adapt Fisher's principles (1935), collect metadata from experiment and sequencing,
- **Rule 3** : Choose a priori tools/methods for bioinformatics and statistical analyses,
- **Rule 4** : Draw conclusions on results.

Experimental design

Find genes that are differentially expressed between a normal skin and a damaged skin on mouse

Sample	Condition	RNA extraction date
S1	control	July 12th, 2016
S2	control	July 12th, 2016
S3	control	July 12th, 2016
S4	wound	July 20th, 2016
S5	wound	July 20th, 2016
S6	wound	July 20th, 2016

Confusion between skin status and RNA extraction date :
comparing healthy and damaged skin is comparing RNAs extracted
July 12th and 20th

Experimental design

Find genes that are differentially expressed between a normal skin and a damaged skin on mouse

Sample	Condition	RNA extraction date
S1	control	July 12th, 2016
S2	control	July 20th, 2016
S3	control	July 25th, 2016
S4	wound	July 12th, 2016
S5	wound	July 20th, 2016
S6	wound	July 25th, 2016

One solution : the day effect is evenly distributed across conditions.

More biological replicates or increasing sequencing depth ?

It depends ! [Haas et al., 2012], [Liu et al., 2014]

- DE transcript detection : (+) biological replicates
- Construction and annotation of transcriptome : (+) depth and (+) sampling conditions
- Transcriptomic variants search : (+) biological replicates and (+) depth

Support

- An experimental design using [multiplexing](#),
- Tools for experimental design decisions : Scotty [Busby et al., 2013], RNAseqPower [Hart et al., 2013], PROPER [Wu et al., 2015]

And do not forget : budget also includes cost of biological data acquisition, sequencing data backup, bioinformatics and statistical analysis.

For a good (nice) experiment design ...

Before the experiment

- Ask a precise and well defined biological question
- List all possible biological confounding effects (sex, age, ...)
- Collect samples while taking care of the distribution of unwanted sources of variation across samples
- Include at least three biological replicates per condition. Technical replicates are not necessary
- Distribute samples on lanes and flow cells ...
 - according to the comparisons to be made
 - without introducing a confusion between technical effects and the biological effects of interest
 - applying the same multiplexing rate on all samples

SARtools

SARTools : Statistical Analysis of RNA-Seq Tools [Varet et al., 2016]

- exports the results into easily readable **tab-delimited files**
 - generates a **HTML report** which displays all the figures produced, explains the statistical methods and gives the results of the differential analysis.
-
- Exploratory data analysis
 - Differential analysis including normalization and multiple testing

Available on R and Galaxy

Exploratory data analysis

Sample comparison for RNA-Seq [Schulze et al., 2012]

Pearson's correlation coefficient

- widely used ...
- ...but highly dependent on sequencing depth and the range of expression samples inherent to the sample.

SERE : Simple Error Ratio Estimate

- ratio of observed variation to what would be expected from an ideal Poisson experiment
- interpretation unambiguous regardless of the total read count or the range of expression
- score of 1 : faithful replication
- score of 0 : data duplication
- scores > 1 true global differences between RNA-Seq libraries

Exploratory data analysis

Multivariate exploratory data analysis

Main goal : explore the structure of the dataset to better understand the proximity between samples and detect possible problems. **This is a quality control step**

Two main tools

- Principal Component Analysis (PCA) or MultiDimensional Scaling (MDS)
- Clustering

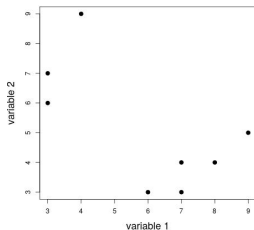
Dimensionality reduction

Problem : n individuals, p genes

$$X = \begin{bmatrix} X_{11} & \dots & X_{1n} \\ X_{21} & \dots & X_{2n} \\ \dots & \dots & \dots \\ X_{p1} & \dots & X_{pn} \end{bmatrix}$$

x_{ij} : value of variable j
for individual i .

Possibility to visualize pair-wise relations by
scatter plots :



When p is large, this is not efficient !

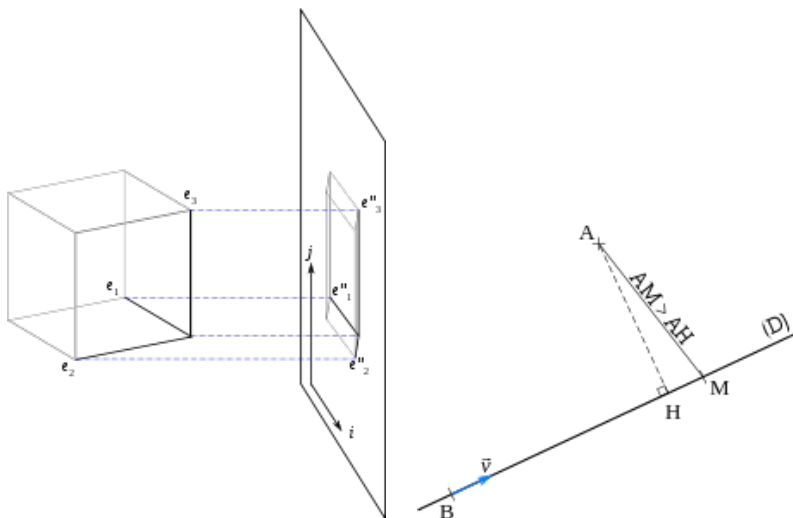
Principal components analysis

Principal component analysis (PCA) :

Main goal : explore the structure of the dataset to better understand the proximity between samples and detect possible problems → often used as a quality control step

- synthesize information and visualize points in a space of reduced dimension
- describe links between variables and which ones explain most variability
- highlight homogeneous subgroups
- detect aberrant individuals

Analyse en composantes principales



Principal components analysis

Principle :

Find axes on which one can project points to obtain a space of reduced dimension comprehensible by the eye.

Projection is a distorting operation \Rightarrow we begin by looking for an axis on which the cloud of points is distorting the less possible during the projection.

PCA uses a **criterion based on variance** to build new axes, also called **components**, in order to preserve variability.

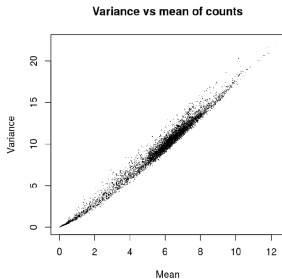
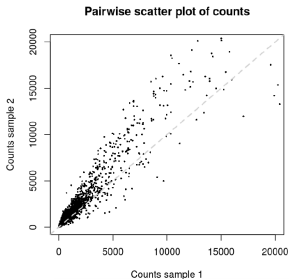
A pre-requisite to apply PCA is to make the data homoscedastic : **the variance must be independent of the intensity.**

Exploratory data analysis

Transformations proposed :

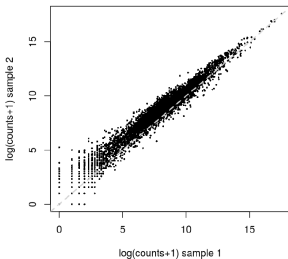
- DESeq2 : VST (Variance Stabilizing Transformation) or rlog (Regularized Log Transformation)
- edgeR : transformation of the count data as moderated log-counts-per-million

Illustration : Without transformation : variance increases with mean

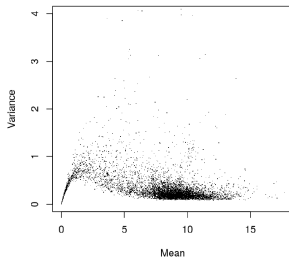


Exploratory data analysis - VST transformation

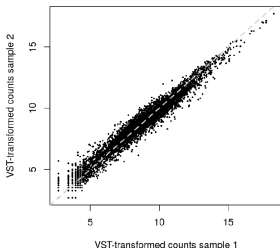
Pairwise scatter plot of log-transformed counts



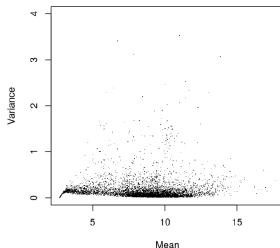
Variance vs mean of log(counts+1)



Pairwise scatter plot of VST-counts



Variance vs mean of VST-transformed counts



Normalization

Definition

Normalization is a process designed to identify and correct **technical biases** removing the least possible biological signal. This step is technology and platform-dependant.

Within-sample normalization

Normalization enabling comparisons of fragments (genes) from a same sample.

No need in a differential analysis context.

Between-sample normalization

Normalization enabling comparisons of fragments (genes) from different samples.

Sources of variability

Read counts are proportional to expression level, gene length and sequencing depth (same RNAs in equal proportions).

Within-sample

- Gene length
- Sequence composition (GC content)

Between-sample

- Depth (total number of sequenced and mapped reads)
- Sampling bias in library construction ?
- Presence of majority fragments
- Sequence composition due to PCR-amplification step in library preparation [Pickrell et al., 2010], [Risso et al., 2011]

Comparison of normalization methods

A lot of different normalization methods...

- Some are part of models for DE, others are 'stand-alone'
- They do not rely on similar hypotheses
- But all of them claim to remove technical bias associated with RNA-seq data

Which one is the best ?

[Dillies et al., 2013], on behalf of StatOmique Group
Evaluation of normalization methods for RNA-Seq differential analysis at the gene level

Comparison of normalization methods

Focus on methods which aim at making read counts comparable across samples

Two main types

- 1 Methods that make read count distributions similar (if not equal)
- 2 Methods assuming that most genes are not differentially expressed

Note that :

- These methods apply on raw (integer) count data, to RNA-seq data (metagenomics), for differential expression analysis
- Other more complex methods have been proposed after the comparison [Risso et al., 2014]
- **Library size** : Number of reads that have been sequenced, mapped and counted for a given sample (sum on columns on the count table)

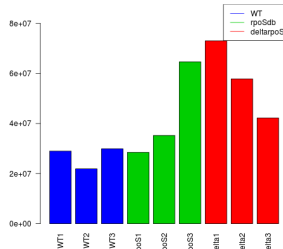


Total Count normalization (TC) [Dudoit et al., 2010]

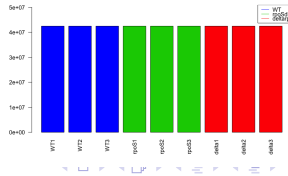
Corrects for differences in the total number of reads

- **Hypothesis** Read count is proportional to gene expression level and sequencing depth (same RNAs with same concentration)
- **But** Very sensitive to the presence of high count genes

SLX, Total Read Count per Sample

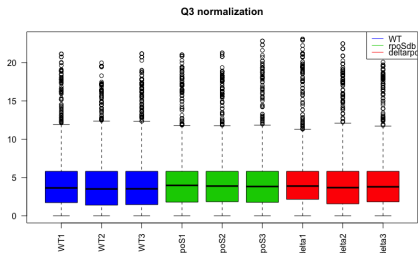


Total Count Normalization, Total Read Count per Sample



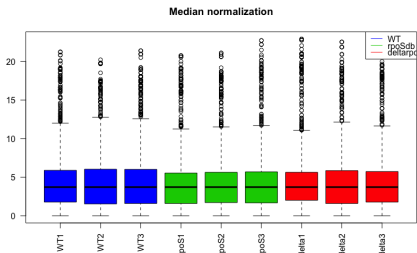
Two variants of Total Counts normalization

Q3 normalization



Third quartile (Q3) is equal across samples

Median normalization



Median is equal across samples

RPKM Normalization [Mortazavi et al., 2008]

Reads Per Kilobase per Million mapped reads

- **Hypothesis** read counts are proportional to gene expression level, gene length and sequencing depth (same RNAs in equal proportions)
- **Method** divide gene count by total count (in million reads) and gene length (in kilobases)
- allows comparisons of gene expression levels within samples

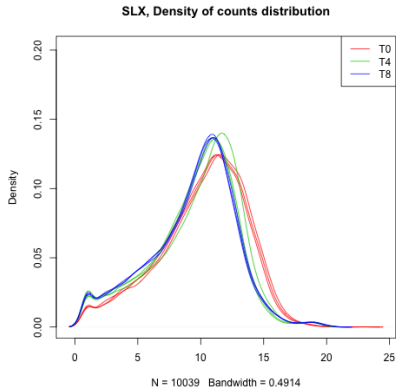
Variations (not compared in [Dillies et al., 2013]) : FPKM, TPM

<http://www.rna-seqblog.com/>

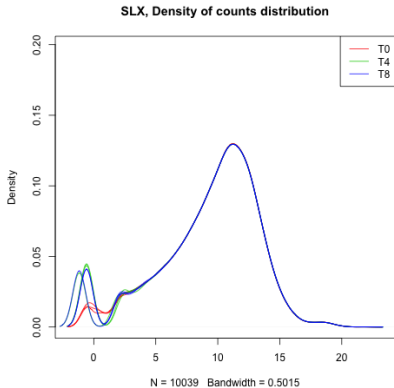
[rpkm-fpkm-and-tpm-clearly-explained/](http://www.rna-seqblog.com/rpkm-fpkm-and-tpm-clearly-explained/)

(Full) Quantile Normalization (FQ)

Hypothesis Read counts have similar distributions across samples



Raw data



Normalized data

Full quantile normalization

Exercise : Calculate the means row by row and for each column, replace the values by the rank in the column.

3	9	3	8	7	7	6	4
6	5	7	4	3	4	3	9
9	4	8	3	8	8	9	6
4	8	4	5	4	9	5	7
7	6	6	7	9	3	4	8

Normalization

5.1	6.9	5.1	6.9	5.9	5.9	6.2	5.1
5.9	5.8	6.2	5.8	5.1	5.8	5.1	6.9
6.9	5.1	6.9	5.1	6.2	6.2	6.9	5.8
5.8	6.2	5.8	5.9	5.8	6.9	5.9	5.9
6.2	5.9	5.9	6.2	6.9	5.1	5.8	6.2

Caution : this normalization provides very good boxplots but can heavily change the measures. It can also favor null variances on rows. Be careful when using it, if not recommended by the platform which generated the data.

"Effective Library Size" [Robinson and Oshlack, 2010]

Motivation

Different biological conditions may express different RNA repertoires, associated with different quantities of total RNAs

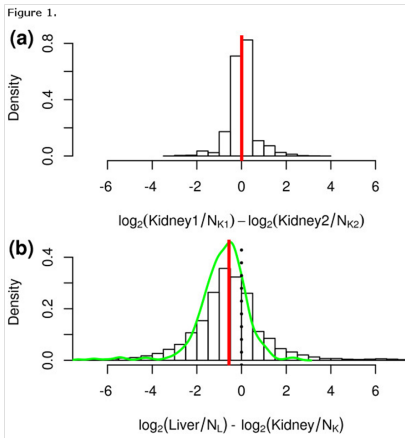
Hypothesis

Most genes are constant across biological conditions

Trimmed Mean of M values (TMM)

[Robinson and Oshlack, 2010]

- Log ratios should be distributed around 0
- Filter on transcripts with null counts, on the resp. 30% and 5% more extreme M_i and A_i
- calculate scaling factors to normalize library sizes



Robinson MD and A. Oshlack. *Genome Biology* 2010, 11:R25.

$$M_i = \log_2\left(\frac{x_{ik1}/N_{k1}}{x_{ik2}/N_{k2}}\right) \text{ and } A_i = 0.5 \times [\log_2(x_{ik1}/N_{k1}) + \log_2(x_{ik2}/N_{k2})]$$

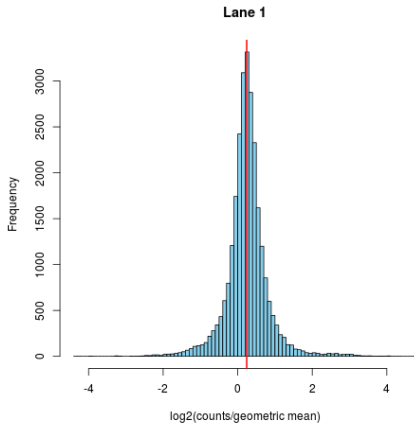
DESeq Normalization [Anders and Huber, 2010]

Normalization factor (for read counts) computed upon genes with a non zero read count in at least one condition

$$\hat{s}_j = \text{median}_i \left(\frac{x_{ij}}{(\prod_{v=1}^n x_{iv})^{1/n}} \right)$$

x_{ij} number of reads in sample j assigned to gene i ,

n number of samples in the experiment



Normalization summary

Methods that compute a normalization factor per sample

Notation

- x_{ij} : number of reads for gene i in sample j
- N_j : number of reads in sample j (library size of sample j)
- n : total number of samples
- \hat{s}_j : normalization factor for sample j
- \hat{x}_{ij} : normalized read count
- \hat{f}_j : scaling factor computed by TMM
- N'_j : library size of sample j normalized with TMM

Normalization summary

Methods that compute a normalization factor per sample

Total count	TMM
$\hat{s}_j = \frac{N_j}{\frac{1}{n} \sum_l N_l}$	$N'_j = N_j * \hat{f}_j, \quad \hat{s}_j = \frac{N'_j}{\frac{1}{n} \sum_j N'_j}$
Q3	DESeq
$\hat{s}_j = \frac{Q3_j}{\frac{1}{n} \sum_l Q3_l}$	$\hat{s}_j = \text{median}_i \frac{x_{ij}}{(\prod_{\nu=1}^n x_{i\nu})^{1/n}}$
Median	Computing normalized reads
$\hat{s}_j = \frac{\text{med}_j}{\frac{1}{n} \sum_l \text{med}_l}$	$\hat{x}_{ij} = \frac{x_{ij}}{\hat{s}_j}$

Which method should I use? [Dillies et al., 2013]

In most cases

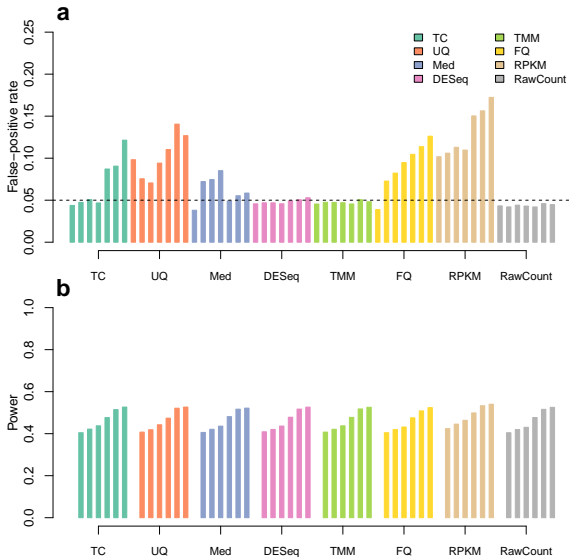
All methods provide comparable results

Anyway ...

Clear differences appear in the presence of high count genes or when the expressed RNA repertoire varies notably across samples

Method	Distribution	Intra-Variance	Housekeeping	Clustering	False-positive rate
TC	-	+	+	-	-
UQ	++	++	+	++	-
Med	++	++	-	++	-
DESeq	++	++	++	++	++
TMM	++	++	++	++	++
FQ	++	-	+	++	-
RPKM	-	+	+	-	-

Which method should I use? [Dillies et al., 2013]



Plan

- 1 Experimental design
- 2 Exploratory data analysis
- 3 Normalization
- 4 Differential analysis**
- 5 Multiple testing
- 6 Gene Set Enrichment Analysis

Statistical significance and practical importance

Differential analysis :

Detect differentially expressed genes between two conditions

Fold change : measure describing how much a quantity changes. Various definitions (see Wikipedia, ipfs.io). In this course : ratio between measurements. If condition A measures 50 and condition B measures 100, fold change = $100/50 = 2$ and measure B is twice higher than measure A.

Log fold change : mean of normalised values in condition 1 - mean of normalised values in condition 2 ($\log B/A = \log B - \log A$)

Question : Why not only using the fold change or log fold change to find differentially expressed genes ?

Statistical significance and practical importance

Fold change does not take the variance of the samples into account. Problematic since variability in omic data is partially marker-specific.

The difference between 102 and 100 is the same as between 4 and 2 but does not seem to have the same importance, regarding the baseline value.

Practical importance and statistical significance have little to do with each other.

- An effect can be important, but undetectable (statistically insignificant) because the data are few, irrelevant, or of poor quality.
- An effect can be statistically significant (detectable) even if it is small and unimportant, if the data are many and of high quality.

Differential analysis

Aim : Detect differentially expressed genes between two conditions

- Discrete quantitative data
- Few replicates
- Overdispersion problem

Challenge : method which takes into account overdispersion and a small number of replicates

- Proposed methods : edgeR, DESeq for the most used and known [Anders et al., 2013]
- An abundant litterature
- Comparison of methods : [Pachter, 2011], [Kvam and Liu, 2012], [Soneson and Delorenzi, 2013], [Rapaport et al., 2013]

Hypothesis testing

Definition

A general method for testing a claim or hypothesis about a parameter in a population, using data measured in a sample.

Four ingredients

- 1 Experimental **data** x_1, x_2, \dots, x_n
- 2 **Statistical model** : assumptions about the independence or distributions of the observations with parameter θ
- 3 **Hypothesis** to test : assumption about one parameter of the distribution
- 4 **Region of rejection** (or critical region) : the set of values of the test statistic T for which the null hypothesis H_0 is rejected. $T = f(X_1, X_2, \dots, X_n)$ is a function which summarizes the data without any loss of information about θ . The distribution of T under H_0 is known.

Critical region and p-value

p-value $p(t)$

For a realisation t of the T test statistic $p(t)$ is the probability (calculating under H_0) of obtaining a test statistic at least as extreme as the one that was actually observed.

In bilateral case :

$$p(t) = \mathbb{P}_{H_0}\{|T| \geq |t|\}$$

The p-value measures the agreement between H_0 and obtained result.

[Link with the critical region](#)

$$\mathbb{P}_{H_0}\{T \in \mathcal{R}\} = \mathbb{P}\{p(t) \leq \alpha\}$$

with α the significance level.

Differential analysis gene-by-gene- with replicates

For each gene i

Is there a significant difference in expression between condition A and B?

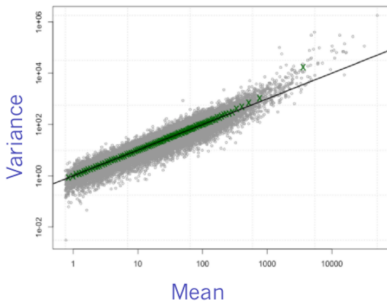
- Statistical model (definition and parameter estimation) - Generalized linear framework
- Hypothesis to test : H_{0i} ; Equality of relative abundance of gene i in condition A and B vs H_{1i} ; non-equality
- Critical region - Wald Test or Likelihood Ratio Test

The Poisson distribution to model counts

- Discrete probability distribution used to describe the number of occurrences of rare events during a given time interval
- Property : Mean = Variance

Mean-Variance Relationship

Technical replicates

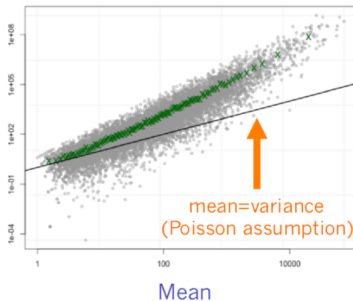


data from Marioni et al. *Gen Res* 2008

From D. Robinson and D. McCarthy

From D. Robinson and D. McCarthy

Biological replicates



data from Parikh et al. *Genome Bio* 2010

Overdispersion in RNA-seq data

Counts from biological replicates tend to have variance exceeding the mean (= overdispersion). Poisson describes only technical variation.

What causes this overdispersion ?

- Correlated gene counts
- Clustering of subjects
- Within-group heterogeneity
- Within-group variation in transcription levels
- Different types of noise present...

In case of overdispersion, increase of the type I error rate (probability to declare incorrectly a gene DE).

Types of noise in data

- 1 Shot noise : unavoidable noise inherent in counting process (dominant for weakly expressed genes) → *well-modeled by Poisson distribution*
- 2 Technical noise : from sample preparation and sequencing, hopefully negligible
- 3 Biological noise : unaccounted for differences between samples (dominant for strongly expressed genes)

Need of an extra-parameter to model the variance

The Negative Binomial Model

Let be X_{ijk} the count for replicate j in condition k from gene i

- X_{ijk} follows a Negative Binomial ($\mu_{ijk} = M_j * \lambda_{ik}$, σ_{ijk}), with M_j library size and λ_{ik} relative abundance of gene i .
- $\sigma_{ijk} = \mu_{ijk}(1 + \phi_i * \mu_{ijk})$

The Negative Binomial distribution

Bernoulli trial

Random experiment with exactly two possible outcomes : success (S) or failure (F)

p : probability of success

Negative Binomial distribution

Repeat Bernoulli trials with probability p of success. NB describes the distribution of the number of failures k before getting n successes

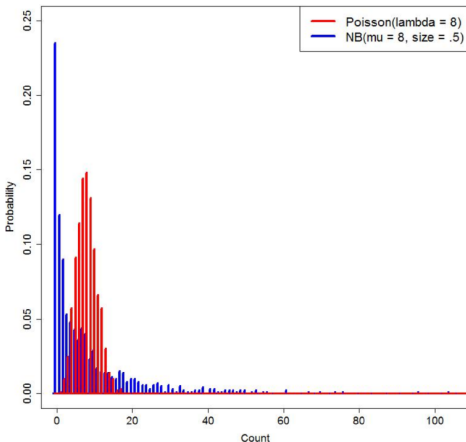
From Poisson to NB

A Negative Binomial distribution is a mixture of Poisson laws with variable parameter. **It is a robust alternative to Poisson in the case of over-dispersed data** (the variance is higher than the mean)

Negative Binomial Models

A supplementary dispersion parameter ϕ to model the variance

Poisson vs Negative Binomial models



Available tests

Models of count data

- Data transformation and gaussian-based model : limma - voom
- Poisson : TSPM
- Negative Binomial : edgeR, DESeq(2), NBPSseq, baySeq, ShrinkSeq, ...

Statistical approaches

- Frequentist Approach : edgeR, DESeq(2), NBPSseq, TSPM, ...
- Bayesian Approach : baySeq, ShrinkSeq, EBSeq, ...
- Non-parametric approach : SAMSeq, NOISeq, ...

Comparison of two conditions

[Soneson and Delorenzi, 2013]

A comparison of methods for differential analysis of RNA-Seq data
[Soneson and Delorenzi, 2013]

- 11 statistical tests included in the study
- R packages
- input data are raw counts (gene-level analysis)
- TMM or DESeq normalization

Main results

- **With only two biological replicates, all the methods show low performances.** They either lack power or poorly control the false positive rate.
- No method outperforms the others in all circumstances : **the method should be chosen according to the dataset**

How to choose ?

- Number of replicates of the experiment
- Presence / absence of outliers
- Constant / variable within-group dispersion
- Balanced / unbalanced differential expression
(results are more accurate and less variable between methods if DE genes are regulated in both directions)
- Simple / complex experiment design

edgeR and DESeq(2)

DESeq2 et edgeR : similarities ...

- Easy to use and well documented R packages
- A 3-step analysis process : normalization, dispersion estimation, statistical test
- Negative Binomial distribution of counts and Generalized Linear Models (GLM) : allows analysis of simple and complex designs

... and differences

- outlier detection and processing
- low counts filtering
- **dispersion estimation**

In both cases, the version matters

Estimating the dispersion : the key question

Coefficient of variation (CV)

Normalized measure of dispersion, ratio of the standard deviation to the mean

In the negative binomial model

$$CV^2 = CV_{\text{technique}}^2 + CV_{\text{biologique}}^2 \quad (1)$$

$$= \frac{1}{\mu_{ijk}} + \phi_i \quad (2)$$

Consequence

Technical variability is the main source of variability in low counts, whereas biological variability is dominant in high counts

Estimating the dispersion : the key question

Problem

Estimate a reliable dispersion from a very small number of replicates (sometimes less than 5)

Why using sophisticated approaches ?

- gene-specific tests \Rightarrow lack of sensitivity (proportion of true positives among positives) due to the lack of information
- common dispersion parameter for all tests \Rightarrow many false positives

Example : empirical bayesian approaches = compromise between gene-specific and common dispersion parameter estimation

Empirical bayesian approaches

Principles

- Bayes theorem : $P(A/B) = P(B/A)P(A)$
- "empirical" \Rightarrow priors from the observed data

$$\tilde{\theta}_g = \hat{\theta}_c + b(\hat{\theta}_g - \hat{\theta}_c)$$

with $\tilde{\theta}_g$ = shrinkage estimator

$\hat{\theta}_c$ = estimator of the mean population

$\hat{\theta}_g$ = usual empirical estimator gene by gene

b = shrinkage factor

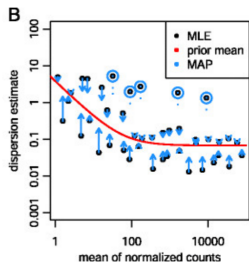
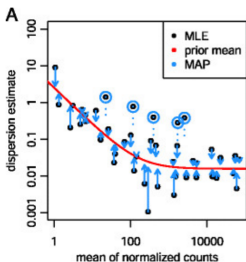
$$b = 1 \Rightarrow \tilde{\theta}_g = \hat{\theta}_g$$

$$b = 0 \Rightarrow \tilde{\theta}_g = \hat{\theta}_c$$

Dispersion estimation with DESeq2

Hypothesis : genes of similar average expression strength have similar dispersion

- 1 Estimate **gene-wise dispersion** estimates using maximum likelihood (ML) (black dots)
- 2 Fit a **smooth curve** (red line)
- 3 **Shrink** the gene-wise dispersion estimates (empirical Bayes approach) toward the values predicted by the curve to obtain final dispersion values (blue arrow heads).



Dispersion estimation with edgeR

- 1 Estimate **gene-wise dispersion** estimates using ML
- 2 Estimate a **common dispersion** parameter by ML
- 3 **Moderate** gene-wise dispersion estimates toward a common estimate or toward a local estimate from genes with similar expression strength using a weighted conditional likelihood.

Differences :

- DESeq2 estimates the width of the prior distribution from the data and therefore automatically controls the amount of shrinkage based on the observed properties of the data.
- edgeR requires a user-adjustable parameter, the prior degrees of freedom, which weights the contribution of the individual gene estimate and edgeR's dispersion fit.

Differences between edgeR and DESeq(2)

Method	Variance	Reference
DESeq, DESeq2	$\mu(1 + \phi\mu)$	[Anders and Huber, 2010], [Love et al., 2014]
edgeR	$\mu(1 + \phi\mu)$	[Robinson et al., 2009]

- **edgeR** : borrow information across genes for stable estimates of ϕ
3 ways to estimate ϕ (common, trend, moderated)
- **DESeq** : data-driven relationship of variance and mean estimated using parametric or local regression for robust fit across genes
- **DESeq2** : relationship of variance and mean (as in DESeq) + dispersion and fold change shrinkage (for PCA and Gene Set Enrichment Analysis) + detection of outliers

DESeq will stop being maintained in a near future, use DESeq2 instead

Robustness - edgeR and DESeq(2)

- **edgeR** : one option : moderate dispersion less towards trend
Allows dispersions to be driven more by the data
- **DESeq** : take the maximum of the fit (trended) or the feature-specific dispersion
Very robust, but many genes pay a penalty, less powerful.
- **DESeq2** : calculate Cook's distance and filter genes with outliers
Can inadvertently filter interesting genes

Goal [Zhou et al., 2014] : achieve a middle ground between protection against outliers while maintaining high power with observation weights

Robustness - edgeR and DESeq(2)

- Robust edgeR (not by default in R) suffers a tiny bit in power with no outliers, but has good capacity to dampen their effect if present (be careful with reviews which take the value by default of edgeR)
- DESeq's policy on outliers has a global effect, resulting in (sometimes drastic) drop in power
- DESeq2 is very powerful in the absence of outliers, but policy to filter outliers results in loss of power
- edgeR and edgeR robust are a bit liberal (5% FDR might mean 6% or 7%)

Comparison of differential analysis methods

To summarize [Soneson and Delorenzi, 2013]

- Obs 1 : The number of replicates matters! (Differently for different methods)
- Obs 2 : Results are more accurate and less variable between methods if DE genes are regulated in both directions.
- Obs 3 : Outlier counts affect different methods in different ways
- The dispersion estimation method matters!

Comparison of differential analysis methods

[Soneson and Delorenzi, 2013]

- Small number of replicates (2-3) or low expression → be careful !!
- Large number of replicates (10 or so) or very high expression → method choice does not matter much.
- Removing genes with outlier counts or using non-parametric methods reduce the sensitivity to outliers
- Allow tagwise dispersion values
- Normalization methods have problems when all DE genes are regulated in one direction.

Comparison of differential analysis methods

[Rapaport et al., 2013]

Evaluation on methods using SEQC benchmark dataset and ENCODE data.

- Significant differences between methods.
- Array-based methods adapted perform comparably to specific methods.
- Increasing the number of replicates samples significantly improves sensitivity over increased sequencing depth.

Multiple Testing

False positive (FP) : A non differentially expressed (DE) gene which is declared DE.

For all 'genes', we test H_0 (gene i is not DE) vs H_1 (the gene is DE) using a statistical test

Problem

Let assume all the G genes are not DE. Each test is realized at α level

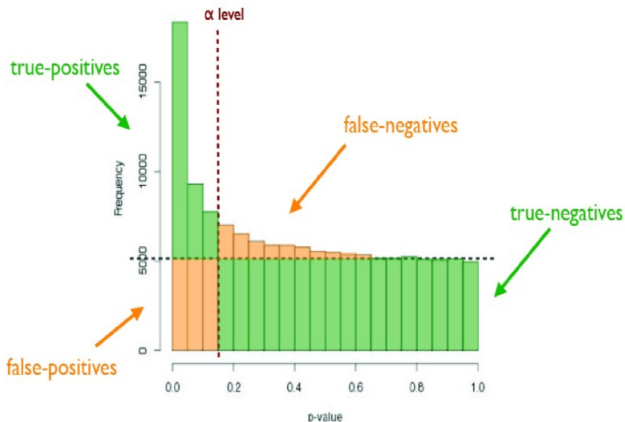
Ex : $G = 10000$ genes and $\alpha = 0.05 \rightarrow E(FP) = 500$ genes.

Simultaneous tests of G null hypotheses

Reality	Declared non diff. exp.	Declared diff. exp.
G_0 non DE genes	True Negatives (TN)	False Positives (FP)
G_1 DE genes	False Negatives (FN)	True Positives (TP)
G Genes	N Negatives	P Positives

Aim : minimize FP and FN .

Standard assumption for p-value distribution



Source : M. Guedj, Pharnext

The Family Wise Error Rate (FWER)

Definition

Probability of having at least one Type I error (false positive), of declaring DE at least one non DE gene.

$$FWER = \mathbb{P}(FP \geq 1)$$

The Bonferroni procedure

Either each test is realized at $\alpha = \alpha^*/G$ level

or use of adjusted pvalue $p_{Bonf_i} = \min(1, p_i * G)$ and $FWER \leq \alpha^*$.

For $G = 2000$ and $\alpha^* = 0.05$; $\alpha = 2.5 \cdot 10^{-5}$.

Easy but conservative and not powerful.

The False Discovery Rate (FDR)

Idea : Do not control the error rate but the proportion of error
⇒ less conservative than control of the FWER.

Definition

The false discovery rate of [Benjamini and Hochberg, 1995] is the expected proportion of Type I errors among the rejected hypotheses

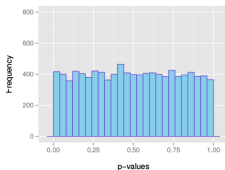
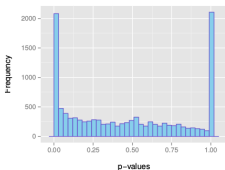
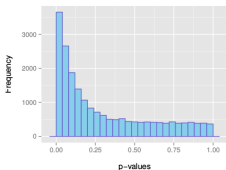
$$\text{FDR} = \mathbb{E}(FP/P) \text{ if } P > 0 \text{ and } 0 \text{ if } P = 0$$

Prop

$$\text{FDR} \leq \text{FWER}$$

p-values histograms for diagnosis

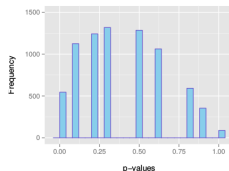
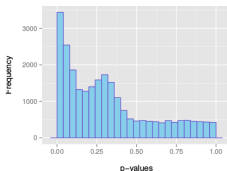
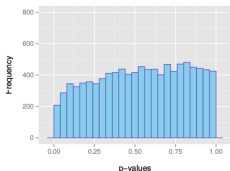
Examples of **expected overall distribution**



- (a) : the most desirable shape
- (b) : very low counts genes usually have large p-values
- (c) : do not expect positive tests after correction

p-values histograms for diagnosis

Examples of not expected overall distribution



- (a) : indicates a batch effect (confounding hidden variables)
- (b) : the test statistics may be inappropriate (due to strong correlation structure for instance)
- (c) : discrete distribution of p-values : unexpected

Multiple testing : key points

- Important to control for multiple tests
- FDR or FWER depends on the cost associated to FN and FP

Controlling the FWER :

Having a great confidence on the DE elements (strong control).

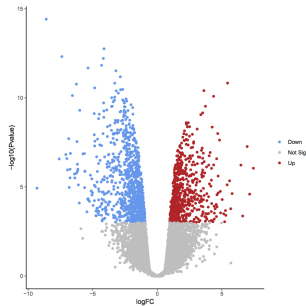
Accepting to not detect some elements (lack of sensitivity \Leftrightarrow a few DE elements)

Controlling the FDR :

Accepting a proportion of FP among DE elements. Very interesting in exploratory study.

Volcano plot

Compromise between statistical significance and importance.
One can adapt the definition of differentially expressed by saying for example "A gene is declared differentially expressed (DE) if the observed difference between two conditions is statistically significant at 5% and the fold change is higher than 2"



Gene Set Enrichment Analysis

Compute overlaps with other gene sets in MSigDB

Use of the hypergeometric distribution which describes the probability of k successes (random draws for which the object drawn has a specified feature) in n draws, *without replacement*, from a finite population of size N that contains exactly K objects with that feature, wherein each draw is either a success or a failure.

The hypergeometric uses the hypergeometric distribution to identify which gene-sets are over-represented in the list of differentially expressed genes. This test is identical to the corresponding one-tailed version of Fisher's exact test.

GSEA history

History of a very cited procedure implemented in the software available on the Broad Institute website :

- first paper : Mootha et al., Nature Genetics, 2004
- Damian and Gorfine published Statistical concerns about the GSEA procedure, Nature Genetics, 2004
- Subramanian et al., PNAS, 2005 : definition of a normalized enrichment score (NES)

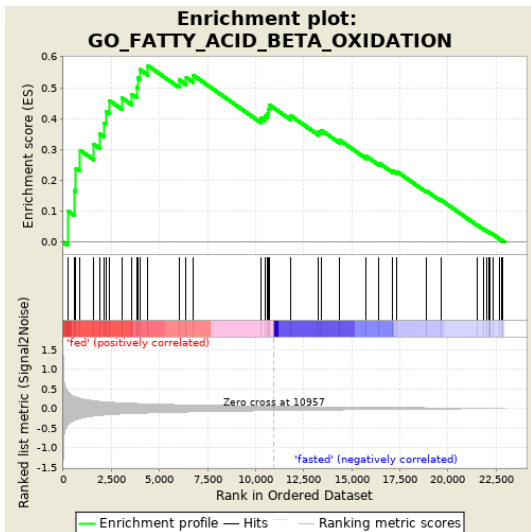
GSEA

The ES reflects the degree to which a set S is over-represented at the extremes (top or bottom) of the entire ranked list L .

The score is calculated by walking down the list L , increasing a running-sum statistic when we encounter a gene in S and decreasing it when we encounter genes not in S .

The magnitude of the increment depends on the ranking metric of the gene with the phenotype. The enrichment score is the maximum deviation from zero encountered in the random walk.

GSEA



GSEA

Estimation of the statistical significance (raw p-value) using phenotype permutations.

- advantage of phenotype permutations : preserving the correlation structure between genes
- not advised to use phenotype permutations when less than 7 samples per condition. In that case, use gene permutations
- in the case of a pre-ranked list, the only possibility is to perform gene permutations

Normalization of the ES for each gene set to account for the size of the set

Adjustment for multiple testing with False Discovery Rate (q-value)

Conclusions

- Methods dedicated to microarrays are not directly applicable to RNA-seq
- Normalization depends on the statistical question
- Include at least 3 replicates per condition for differential analysis
- Large number of replicates (10 or so) or very high expression → method choice of differential analysis does not matter much.
- Removing genes with outlier counts or using non-parametric methods reduce the sensitivity to outliers
- Don't forget to correct for multiple testing!

Conclusions

Adapt the method to your data

Specific methods have been developed for few replicates.

The need for 'sophisticated' methods decreases when the number of replicates increases.

GSEA help find differentially expressed genes when not enough replicates were present in the initial study. Avoid merging the data when a high study effect is expected, prefer an appropriate statistical analysis!

Want to go further ?

To practice more : Galaxy permanences

<https://wikis.univ-lille.fr/bilille/permanences>

To obtain help in statistical analysis of omic data :

bilille call for projects (around december each year, to plan the calendar of engineers)

Scotty : a web tool for designing RNA-Seq experiments to measure differential gene expression.
[Bioinformatics 2013, 29\(5\),656 :657.](#)



[Dillies MA, Rau A, Aubert J, Hennequet-Antier C et al](#)

A comprehensive comparison of normalization methods for Illumina high-throughput RNA-sequencing data analysis
[Briefings in Bioinformatics 2013, 14 :6, 671-683.](#)



[Dudoit S, Maya O and Jacob L.](#)

Short course on RNA seq and ChIP seq data analysis.
Valencia, Nov. 2010.



[Eisenberg EE and Levanon EY.](#)

Human housekeeping genes are compact.
[Trends Genet, 19\(7\) :362-365.](#)



[Fisher RA](#)

The Design of experiments
[Oliver and Boyd 1935, 1-252](#)



[Haas BJ, Chin M, Nusbaum C, Birren BW, Livny J](#)

How deep is deep enough for RNA-Seq profiling of bacterial transcriptomes?
[BMC genomics 2012, 1 \(13\),734.](#)



[Hansen KD, Brenner SE, Dudoit S.](#)

Biases in Illumina transcriptome sequencing caused by random hexamer priming.
[Nucleic Acids Research, 2010, 1-7.](#)



[Hansen KD, Irizarry RA and Wu Z](#)

Removing technical variability in RNA-seq data using Conditional Quantile Normalization
[Biostatistics 2011, 13 :2, pp204-216](#)



[Hart SN, Therneau TM, Zhang Y, Poland GA, Kocher J-P](#)

Calculating Sample Size Estimates for RNA Sequencing Data.
[Journal of Computational Biology 2013, 12\(20\), 970 :978](#)



[Kvam V, Liu P](#)

A comparison of statistical methods for detecting differentially expressed genes from RNA-seq data
[American Journal of Botany 2012 99\(2\), 248-256](#)



[Liu Y, Zhou J, White K](#)

RNA-seq differential expression studies : more sequence or more replication ?
[Bioinformatics 2014, 30\(3\),301 :304.](#)



[Love MI, Huber W and Anders S](#)

Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2.
[Genome Biology 2014, 15 :550.](#)



[Marioni JC, Mason CE et al.](#)

RNA-seq : An assessment of technical reproducibility and comparison with gene expression arrays.
[Genome Research 2008, 18 : 1509-1517](#)



[Marot G, Foulley JL, Mayer CD, Jaffrézic F.](#)

Moderated effect size and P-value combinations for microarray meta-analyses.
[Bioinformatics 2009, 25\(20\) :2692-9.](#)



[McIntyre LM, Lopiano KK, Morse AM, Amin V, Oberg AL, Young LJ, Nuzhdin SV](#)

RNA-seq : technical variability and sampling
[BMC Genomics 2011, 12 :293.](#)



[Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B.](#)

Mapping and quantifying mammalian transcriptomes by RNA-seq.
[Nature Methods, 2008 Jul; 5\(7\) ; 621-628](#)



[Oshlack A and Wakefield MJ](#)

Transcript length bias in RNA-seq confounds systems biology
[Biology Direct 2009, 4 :14.](#)



[Pachter L](#)

Models for transcript quantification from RNA-seq

eprint 2011 arXiv :1104.3889



Pickrell JK, Marioni JC, Pai AA, Degner JF, Engelhardt BE, Nkadori E, Veyrieras JB, Stephens M, Gilad Y, Pritchard JK.

Understanding mechanisms underlying human gene expression variation with RNA sequencing. Nature letters, 2010, vol 464.



Rapaport F, Khanin R, Liang Y, Pirun M, Krek A, Zumbo P, Mason CE, Socci ND, Betel D

Comprehensive evaluation of differential gene expression analysis methods for RNA-seq data Genome Biology 2013 ,14 :R95



Rau A, Marot G, Jaffrézic F

Differential meta-analysis of RNA-seq data from multiple studies. BMC Bioinformatics 2014 ,15 :91



Risso D, Schwartz K, Sherlock G, Dudoit S.

GC-content normalization for RNA-Seq data BMC Bioinformatics 2011, 17, 12 :480



Risso D, Ngai J, Speed T and Dudoit S

Normalization of RNA-seq data using factor analysis of control genes or samples Nature Biotechnology 2014, 32 (9), 896-905



Robinson MD and Smyth, GK.

Moderated statistical tests for assessing differences in tag abundance. Bioinformatics 23(21) ; 2881-2887



Robinson MD and Smyth, GK.

Small-sample estimation of negative binomial dispersion, with applications to SAGE data Biostatistics (2008), 9, 2 ; 321-332



Robinson MD, McCarthy DJ, Smyth, GK.

edgeR : a Bioconductor package for differential expression analysis of digital gene expression data. Bioinformatics 2009



[Robinson MD, Oshlack A.](#)

A scaling normalization method for differential expression analysis of RNA-seq data.
[Genome Biology](#) 2010, 11 :R25



[Robles J.A., Qureshi S.E., Stephen S.J., Wilson S.R., Burden C.J., Taylor J.M.](#)

Efficient experimental design and analysis strategies for the detection of differential expression using RNA-Sequencing
[BMC Genomics](#) 2012, 13 :484



[Schulze SK, Kanwar R, Gölzenleuchter M, Therneau TM, Beutler AS.](#)

SERE : Single-parameter quality control and sample comparison for RNA-Seq
[BMC Genomics](#) 2012, 13 :524



[Soneson C and Delorenzi M](#)

A comparison of methods for differential expression analysis of RNA-seq data
[BMC Bioinformatics](#) 2013, 14 :91



[Su AI, Wiltshire T, Batalov S, Lapp H, Ching KA, Block D, Zhang J, Soden R, Hayakawa M, Kreiman G, Cooke MP, Walker JR, Hogenesch JB.](#)

A gene atlas of the mouse and human protein-encoding transcriptomes.
[Proc. Natl. Acad. Sci. USA](#), 101(16) :6062-6067.



[A. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A, Pomeroy SL, Golub TR, Lander ES, Mesirov JP.](#)

Gene set enrichment analysis : A knowledge-based approach for interpreting genome-wide expression profiles.
[PNAS](#) 2005, 102(43) :15545–15550.



[Trapnell C, Hendrickson D, Sauvageau M, Goff L, Rinn J and Pachter L](#)

Differential analysis of gene regulation at transcript resolution with RNA-seq
[Nature Biotechnology](#) 2013, 31, 1



[Varet H, Brillet-Gueguen L, Coppée JY, Dillies MA](#)

SARTools : A DESeq2- and edgeR-Based R Pipeline for Comprehensive Differential Analysis of RNA-Seq Data

