

Université Paris Sud – Paris Saclay

2018 - 2019

Differential expression analyses: How to (better) understand your results ?

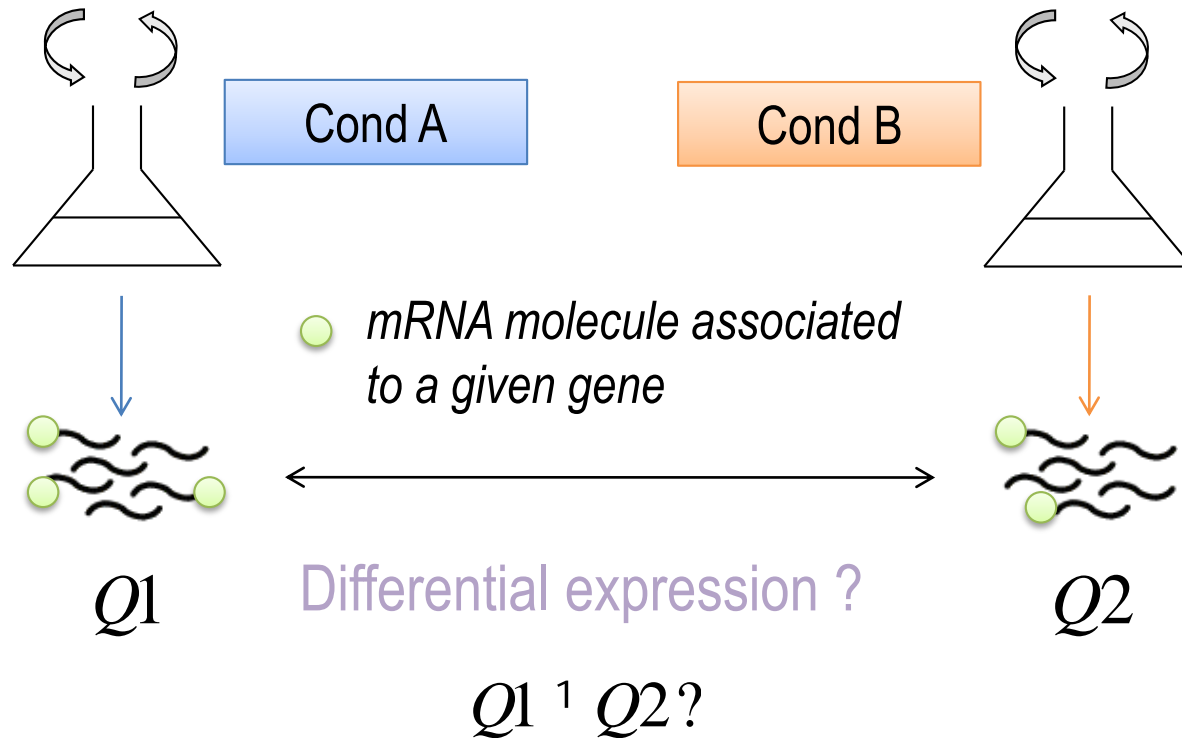
Gaëlle LELANDAIS

gaelle.lelandais@u-psud.fr



This work is licensed under a
Creative Commons Attribution 3.0 Unported License

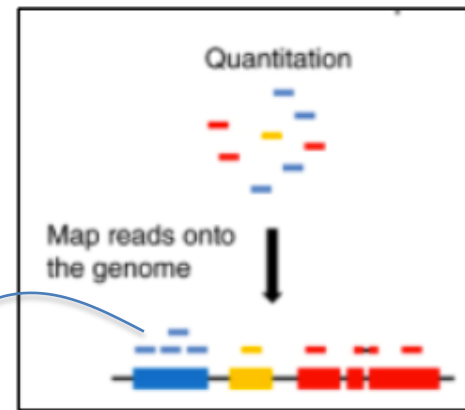
What is the question ?



- All genes are successively analyzed in a differential analysis

Experimental dataset

(number of mapped reads)



Cond A

Cond B

	R1	R2	R3	R1	R2	R3
A1BG	4	6	2	7	6	7
A1CF	41	33	42	32	42	32
A2M	1	3	1	4	3	7
A2ML1	3	2	2	6	7	3
A2MP1	3	2	2	1	1	0
A3GALT2	1	4	4	3	2	1
A4GALT	420	344	291	327	360	371
A4GNT	1	1	2	1	3	3
AA06	0	0	0	0	0	0
AAAS	2452	2192	1977	2054	2134	2100
AACS	3234	2804	2609	1678	1670	1742
AACSP1	1544	1369	1300	1926	2015	1963

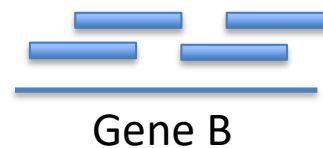
Data normalisation : Why ?

Number of mapped reads is related to gene length

A1BG	4
A1CF	41
A2M	1
A2ML1	3
A2MP1	3
A3GALT2	1
A4GALT	420
A4GNT	1
AA06	0
AAAS	2452
AACS	3234
AACSP1	1544



8 mapped
reads



4 mapped
reads

Data normalisation : Why ?

Number of mapped reads is related to library size

A1BG	4	7
A1CF	41	32
A2M	1	4
A2ML1	3	6
A2MP1	3	1
A3GALT2	1	3
A4GALT	420	327
A4GNT	1	1
AA06	0	0
AAAS	2452	2054
AACS	3234	1678
AACSP1	1544	1926

Sum 1

Sum 2

Systematic bias ?

(If Sum 1 = 40.000 and Sum 2 = 30.000)

Section 1

GENERAL PRINCIPLE

Working hypothesis for normalization

Most of the genes are not
differentially expressed

$$\text{LogFC} \approx 0$$

LogFC (Fold Change) parameter

➤ Average (mean) value of $\log_2(C_A/C_B)$ in 3 replicates

$$\log FC_g = \frac{1}{3} \left(\log_2 \left(\frac{Q1}{Q2} \right)_{g,R1} + \log_2 \left(\frac{Q1}{Q2} \right)_{g,R2} + \log_2 \left(\frac{Q1}{Q2} \right)_{g,R3} \right)$$

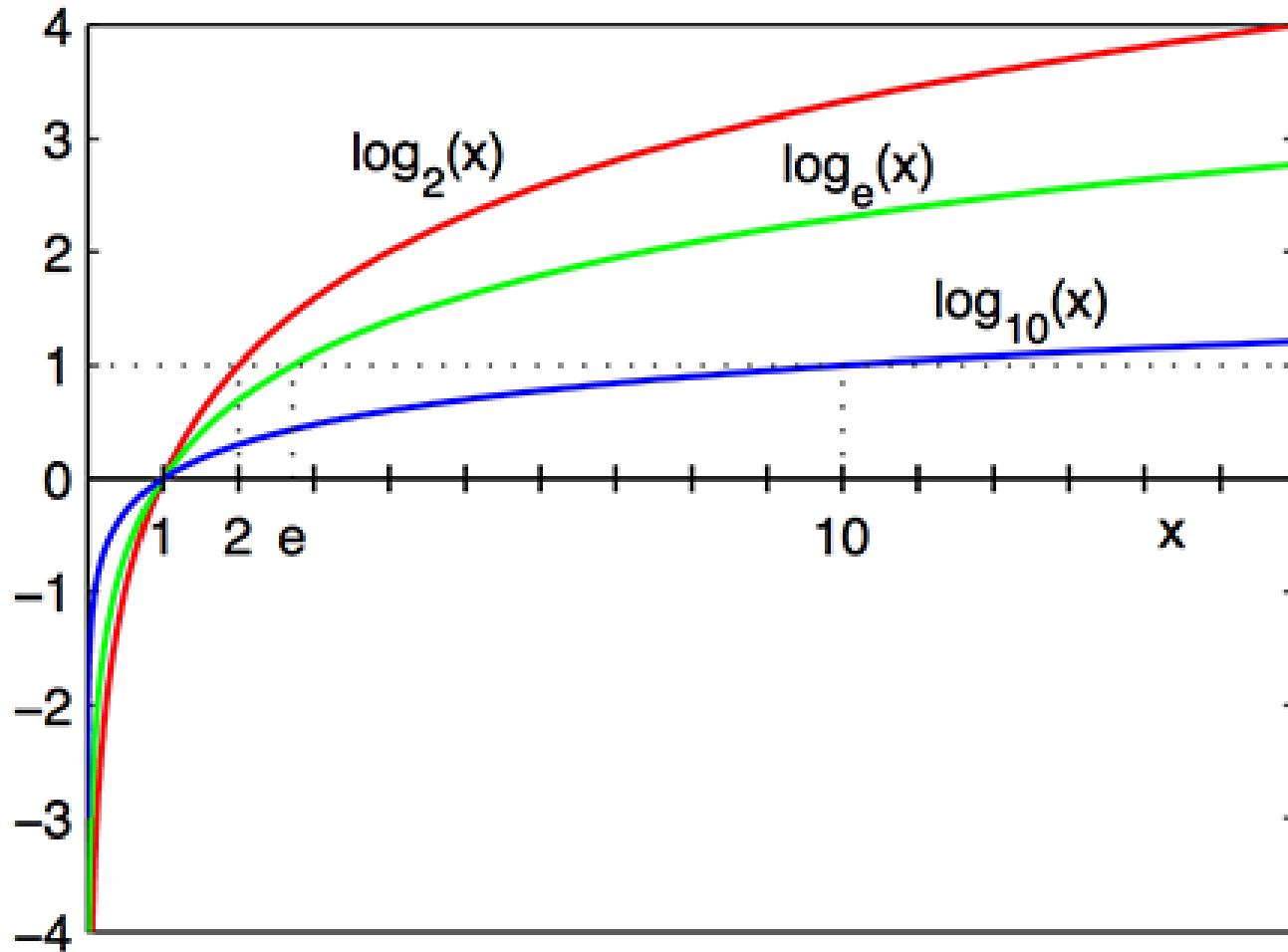
$$\log FC_g > 0 \Leftrightarrow Q1_g > Q2_g$$

 Up-regulated gene

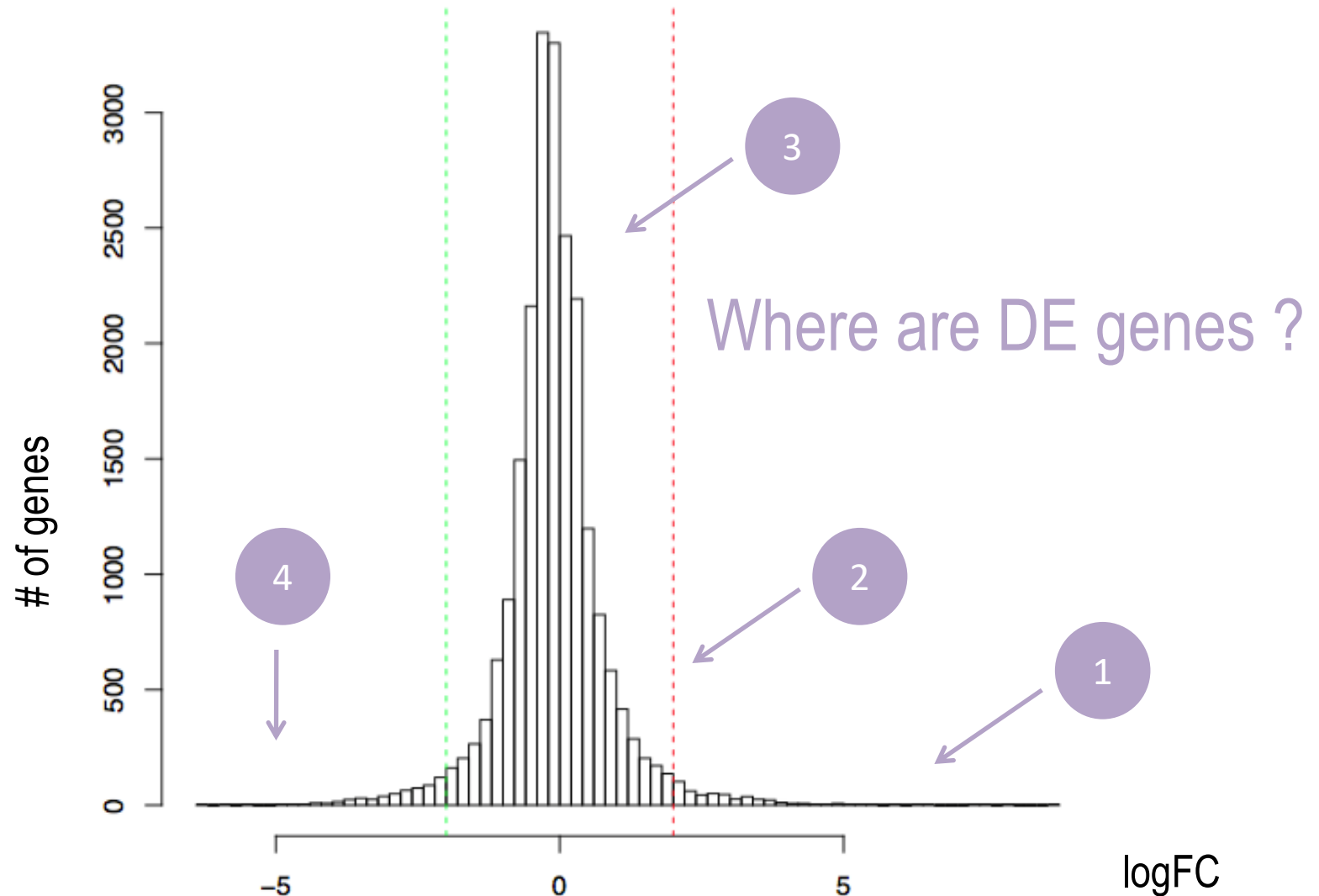
$$\log FC_g < 0 \Leftrightarrow Q1_g < Q2_g$$

 Down-regulated gene

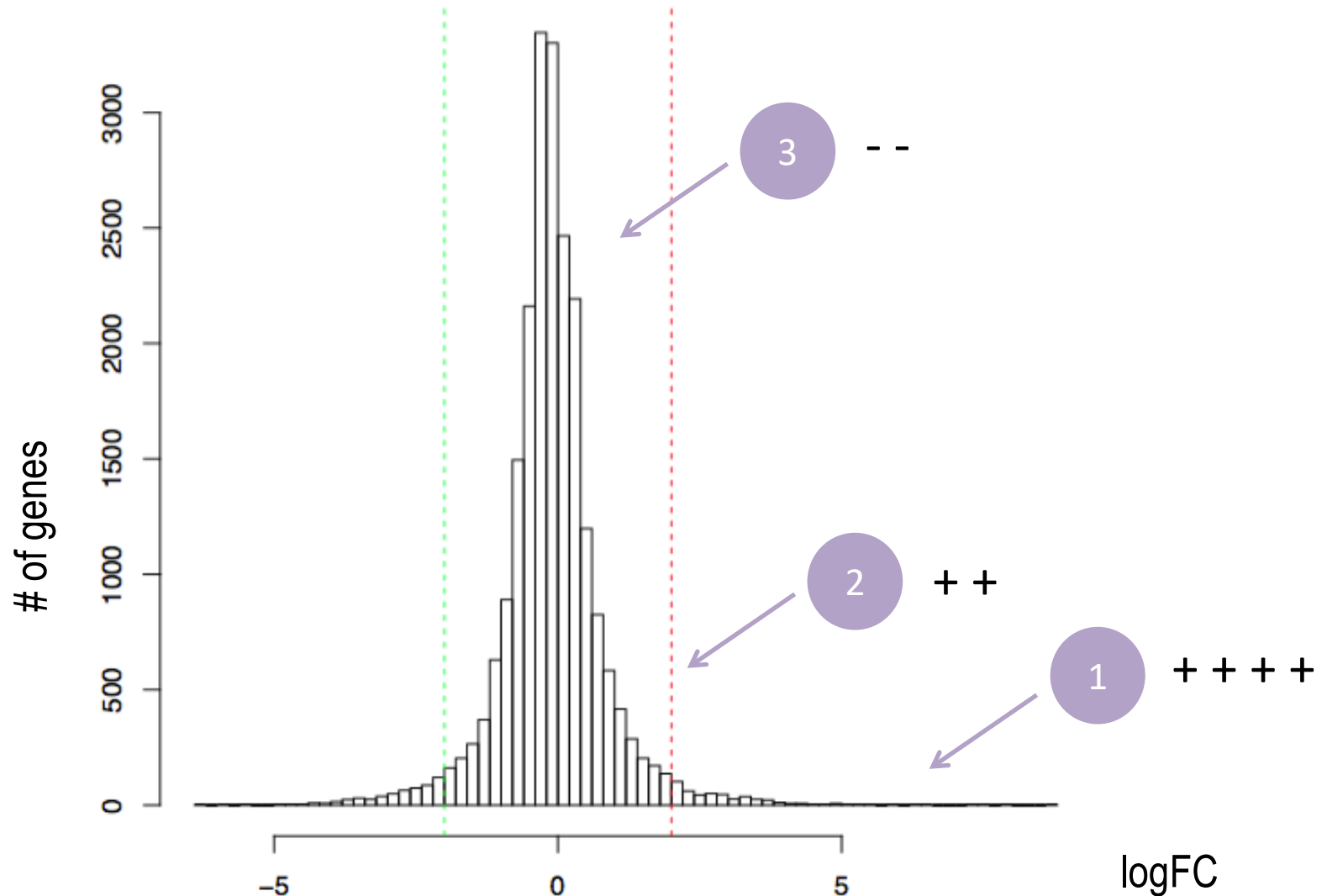
Why using logFC instead of FC ?



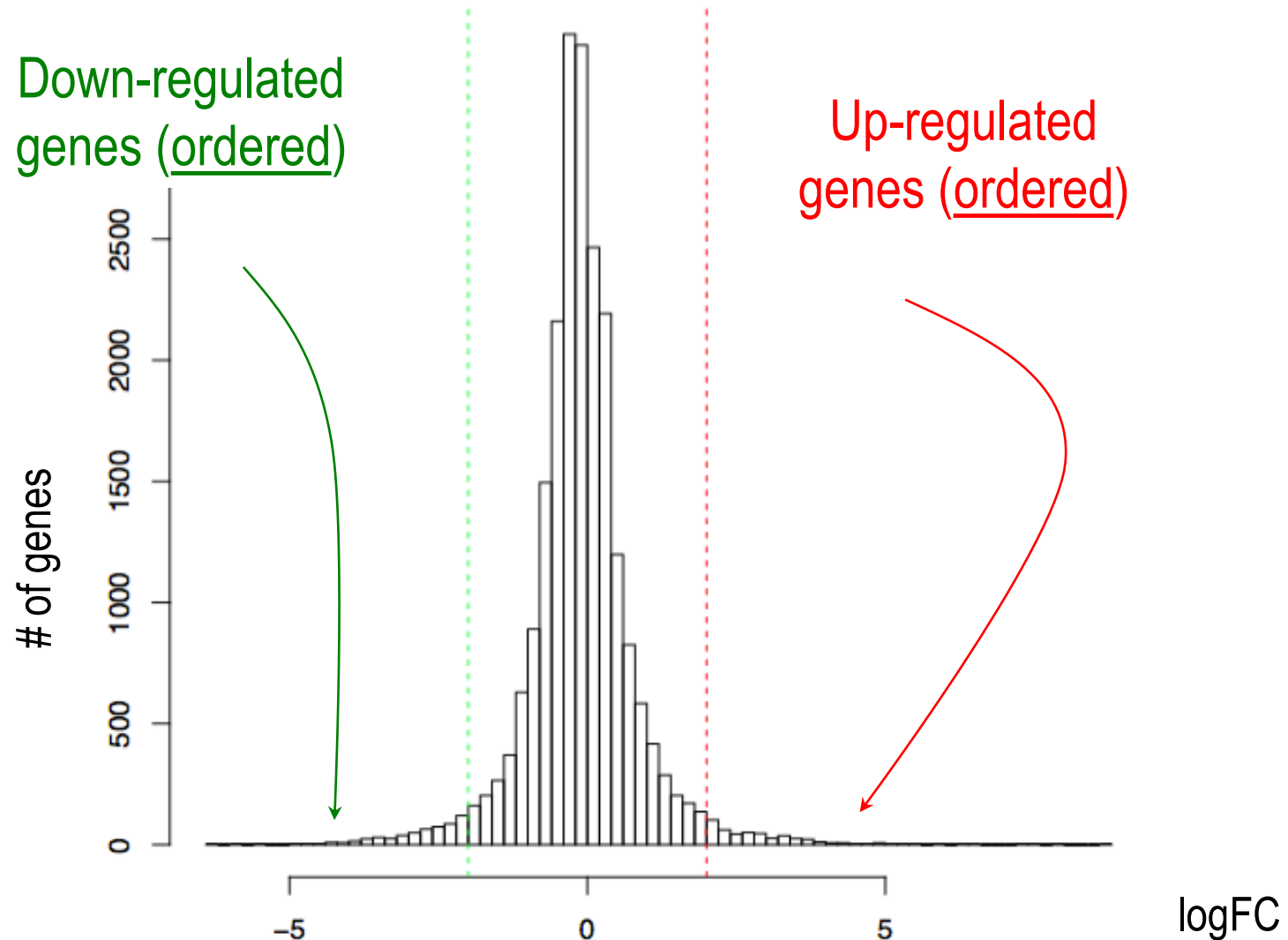
Distribution of logFC values



It's all about “confidence level”



DE analysis based on logFC only



In addition to logFC ...

- Two genes can have identical logFC, but different individual $\log_2(C_A/C_B)$ values

+ + Gene 1: $\log FC_{G1} = \frac{1}{3}(0.5 + 2.5 + 6) = 3$

+ + + + Gene 2: $\log FC_{G2} = \frac{1}{3}(3 + 2.8 + 3.2) = 3$

Data variability

Same level of confidence ?

Quantification of the “data variability”

➤ Variance estimation

➤ Standard deviation

$$Var(X) = \frac{1}{n-1} \sum (x_i - m)^2$$

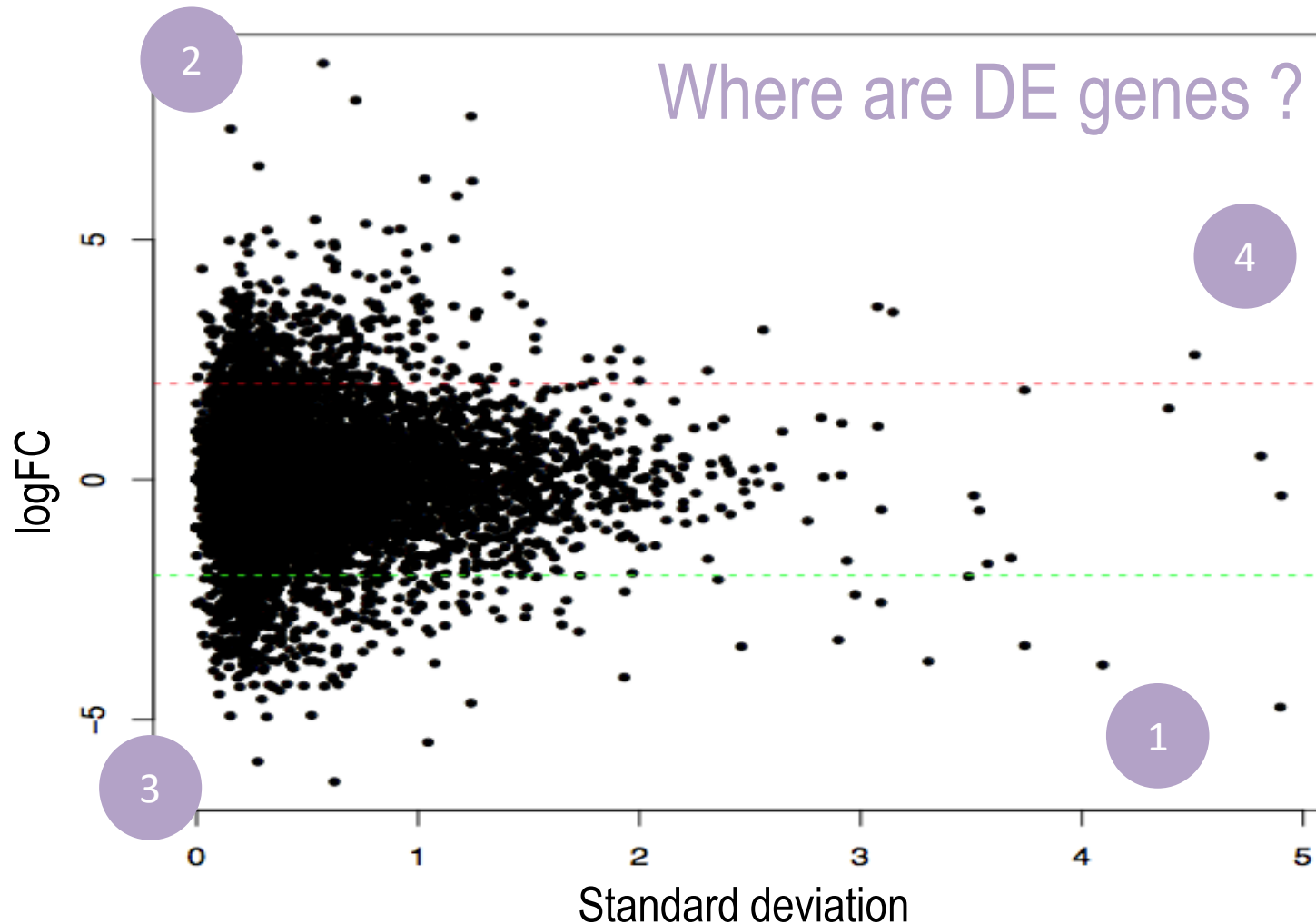
$$SD = \sqrt{Var(X)}$$

Application to logFC :

	logFC	Variance	SD
Gene 1	3	7.75	2.78
Gene 2	3	0.04	0.20



DE analysis based on logFC and SD



Solution to combine differential expression (logFC) and data variability (SD)

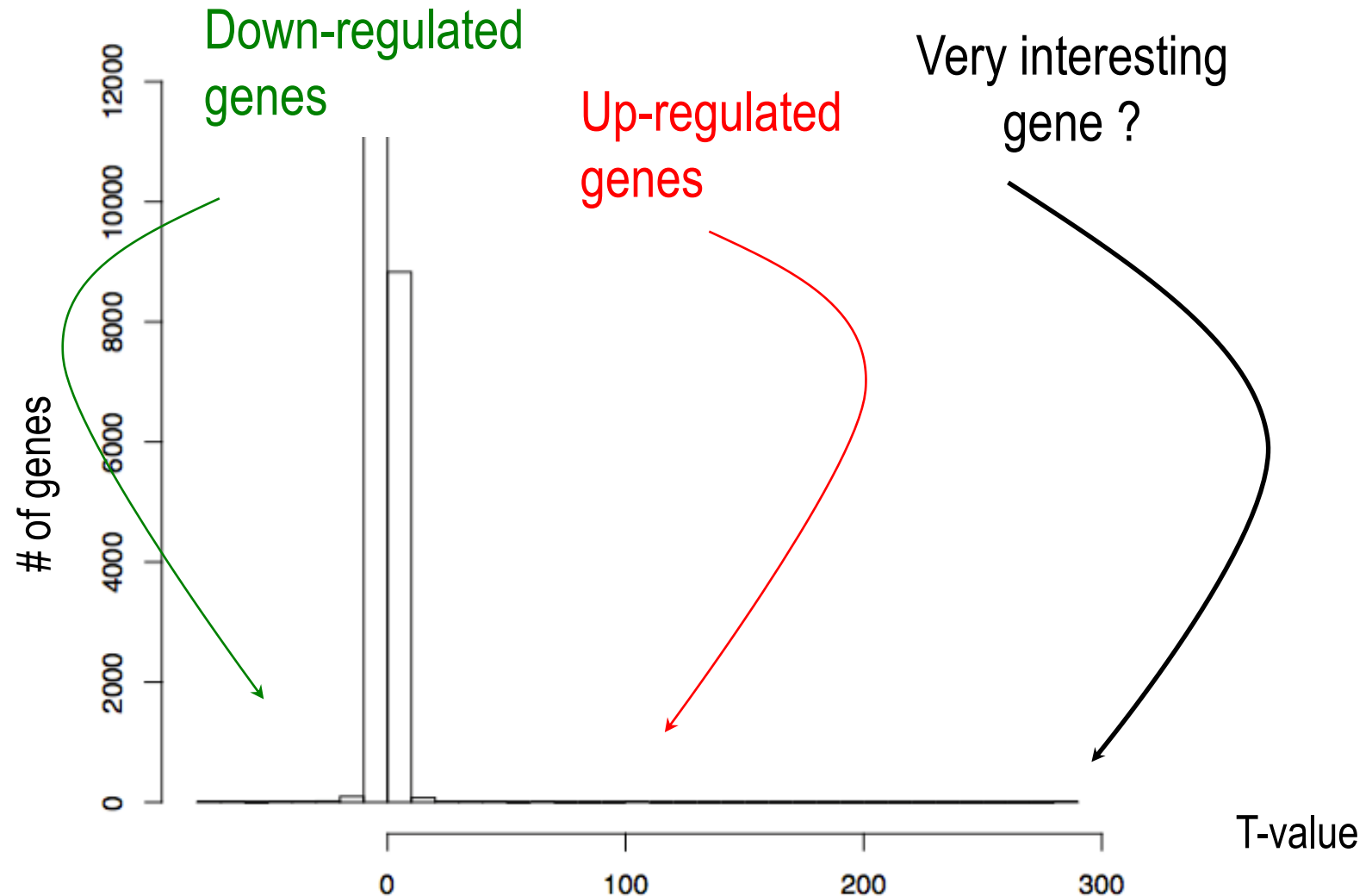
- Student parameter t (used in classical statistics)

$$t_g = \frac{\log FC_g}{\frac{SD_g}{\sqrt{n}}}$$

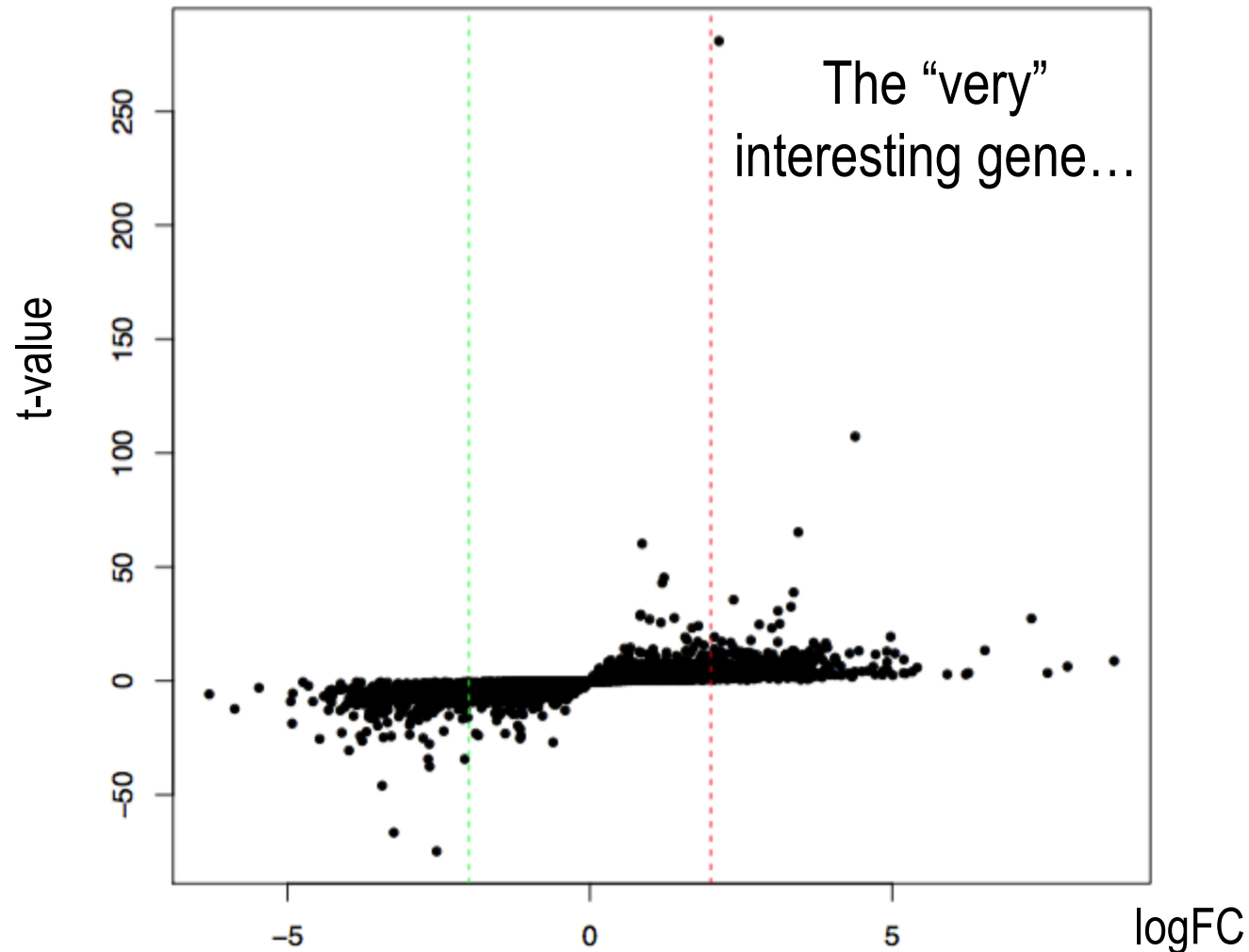
How to obtain high/low values of “t” ?

Does it work correctly ?

Distribution of t-values (logFC)



Drawback associated to the classical t-value

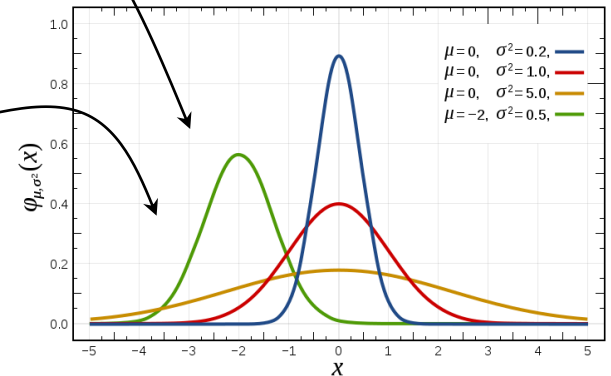


Also, p-value can't be calculated because ...

- The random variable must be normally distributed :

Gene 1: $\log FC_{G1} = \frac{1}{3} (0.5 + 2.5 + 6) = 3$

Gene 2: $\log FC_{G2} = \frac{1}{3} (3 + 2.8 + 3.2) = 3$



This is impossible to test, because of the limited number of replicates

R script to be used/modified

Search_DEgenes.R

```
search_DEgenes_M2Bioinfo.R
search_DEgenes_M2Bioinfo.R No Selection

#####
# Gaelle LELANDAIS <gaelle.lelandais@u-psud.fr>
# Melina GALLOPIN <melina.gallopin@u-psud.fr>
#####

####
# Ce script R est destine aux etudiants de Master. Il a pour objectif de les aider dans la
# comprehension des methodes de selection des genes differentiellement exprimes.
####

####
# Librairies R necessaires : LIMMA, MASS, DESeq2, edgeR
####

#-----
#-----

pdf("graphics.pdf")

# Lecture des données
countData = read.table("count_dataFile.txt", header = T, row.names = 1)
# --> Ces donnees sont en relation avec l'article :
#
# "Determination of a Comprehensive Alternative Splicing Regulatory Network and
# Combinatorial Regulation by Key Factors during the Epithelial-to-Mesenchymal Transition"
#
# Yang et al. (2016, Molecular and Cellular Biology)
#
# --> Deux points de temps ont ete etudies par RNAseq : "Day0" et "Day7", 3 replicats sont
# disponibles pour chaque point de temps.

#-----
# Pretraitement des donnees entre les experiences (filtrage des faibles et normalisation)
#-----

## Suppression des genes sans valeurs (0 pour toutes les experiences)
length(which(rowSums(countData)==0)) ## 2957
countData <- countData[-which(rowSums(countData) == 0),]

# Distribution des valeurs de comptages dans les différents échantillons
boxplot(log(countData + 1), ylab = "#reads (log scale)",
        main = "Boxplot of read counts", col = c(rep("grey", 3), rep("blue", 3)),
        names = c(paste("Day0_", 1:3, sep = ""), paste("Day0_", 1:3, sep = "")))
# --> Les distributions des donnees de comptage sont tres proches entre les differentes
# conditions.
```

SECTION 2

DEDICATED METHODS TO IDENTIFY DIFFERENTIALLY EXPRESSED GENES

R package DESeq2

➤ Read counts $\xrightarrow{\text{DESeq2}}$ ➤ Statistical parameters

RNAseq_counts_parapsilosis.txt

Normoxia_1	Normoxia_2	Normoxia_3	Normoxia_4	Normoxia_5	Normoxia_6
Hypoxia_1	Hypoxia_2	Hypoxia_3	Hypoxia_4		
CPAR2_484750	589	498	630	612	627
CPAR2_212570	559	648	576	689	988
CPAR2_302130	167	83	230	151	180
CPAR2_209760	869	1438	972	1250	1588
CPAR2_483260	2549	2313	3115	2801	5834
CPAR2_209040	533	313	646	532	335
CPAR2_804840	44	42	57	44	68
CPAR2_211500	621	835	736	783	1156
CPAR2_182540	1672	2432	1845	2134	3653
CPAR2_189950	1772	2769	2175	2424	1629
CPAR2_800250	166	144	157	179	181
CPAR2_503900	1023	1435	1320	1414	1894
CPAR2_180740	238	180	279	224	187
CPAR2_206160	705	933	831	823	1129
CPAR2_181420	383	552	448	526	612
CPAR2_805280	378	515	413	500	850
CPAR2_202350	1253	2099	1529	1826	2918
CPAR2_783820	3668	5607	4518	5199	5891
CPAR2_582480	219	291	284	302	233
CPAR2_180810	1118	1908	1332	1616	1293
CPAR2_185300	5288	7896	6291	7105	11283
CPAR2_487740	323	422	422	435	956
CPAR2_802340	640	848	822	900	1686
CPAR2_186320	8875	13180	10310	11985	16549
CPAR2_301210	167	191	218	204	540
CPAR2_302570	743	1121	856	947	1407
CPAR2_583570	99	137	135	101	482

baseMean

DESeq2_statistics.txt

baseMean	log2FoldChange	lfcSE	stat	pvalue	padj
CPAR2_484750	557.116255077326	0.133444946455645	0.190487408943552		
CPAR2_212570	559.374382948797	0.0287350866112832	0.12662336890635		
CPAR2_302130	173.841921340983	0.35023424511046	0.320288521371104		
CPAR2_209760	954.765738326101	-0.114035759306079			
CPAR2_483260	2857.44847832337	0.358640646962554	0.197651989899545		
CPAR2_209040	491.800116175746	0.140998897015593	0.309111624371794		
CPAR2_804840	50.1539296854689	0.345329251780373	0.261316423135524		
CPAR2_211500	604.715757128726	0.0064962161591738	0.119236535927135		
CPAR2_182540	1997.88804218508	0.236796415663859	0.148912116775271		
CPAR2_189950	1772.276921340983	0.1976103167	0.1976103167		
CPAR2_800250	166.144157179	-0.913291	0.0435876423		
CPAR2_503900	1023.14351320	0.0928878	0.1457826686		
CPAR2_783820	3668.56074518	-0.076718	0.84540483756		
CPAR2_582480	219.291284	-0.918803	0.8398147556		
CPAR2_180810	1118.19081332	-0.399161611153727			
CPAR2_185300	5288.78966291				
CPAR2_487740	323.422422				
CPAR2_802340	640.848822				
CPAR2_186320	8875.1318010310				
CPAR2_301210	167.191218				
CPAR2_302570	743.1121856				
CPAR2_583570	99.137135				

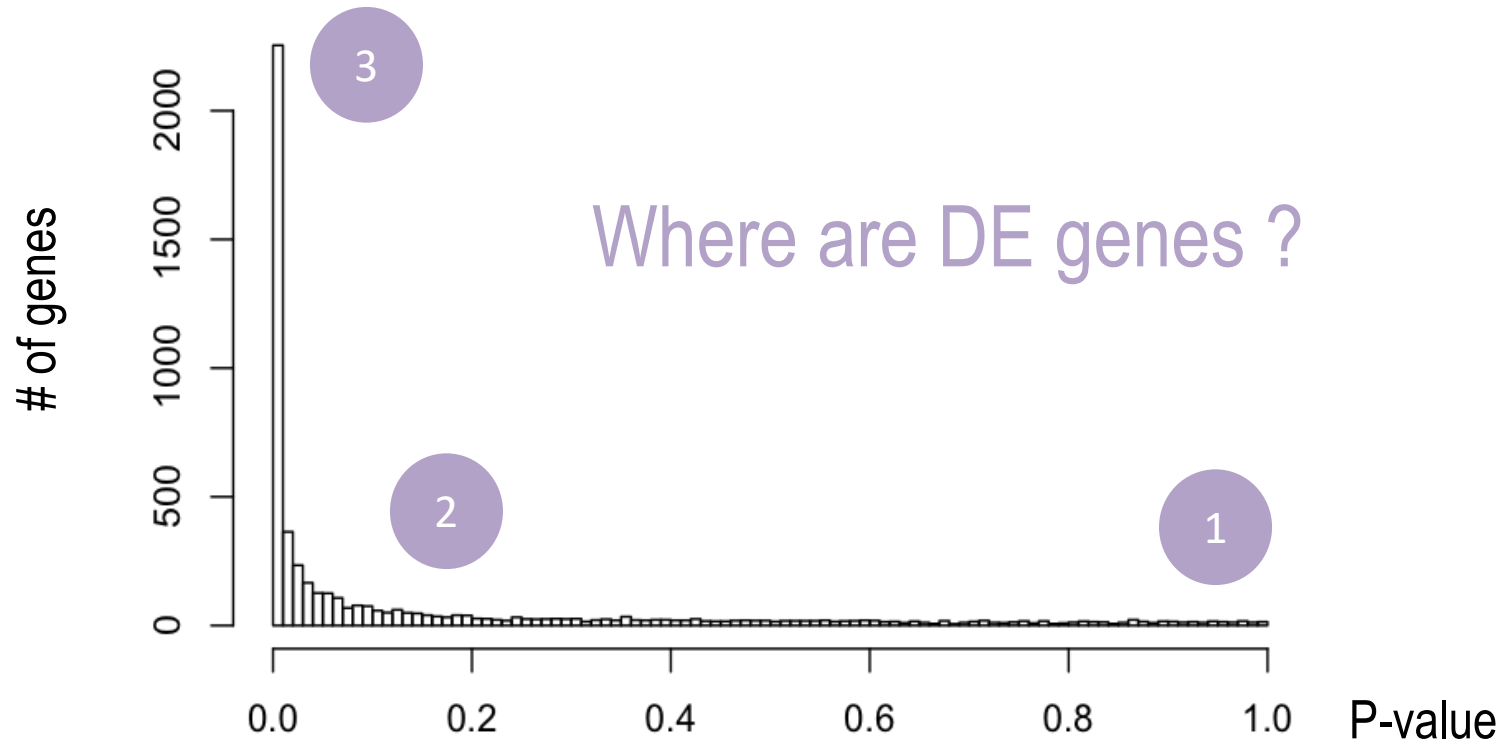
log2Fold
Change

pvalue and
padj



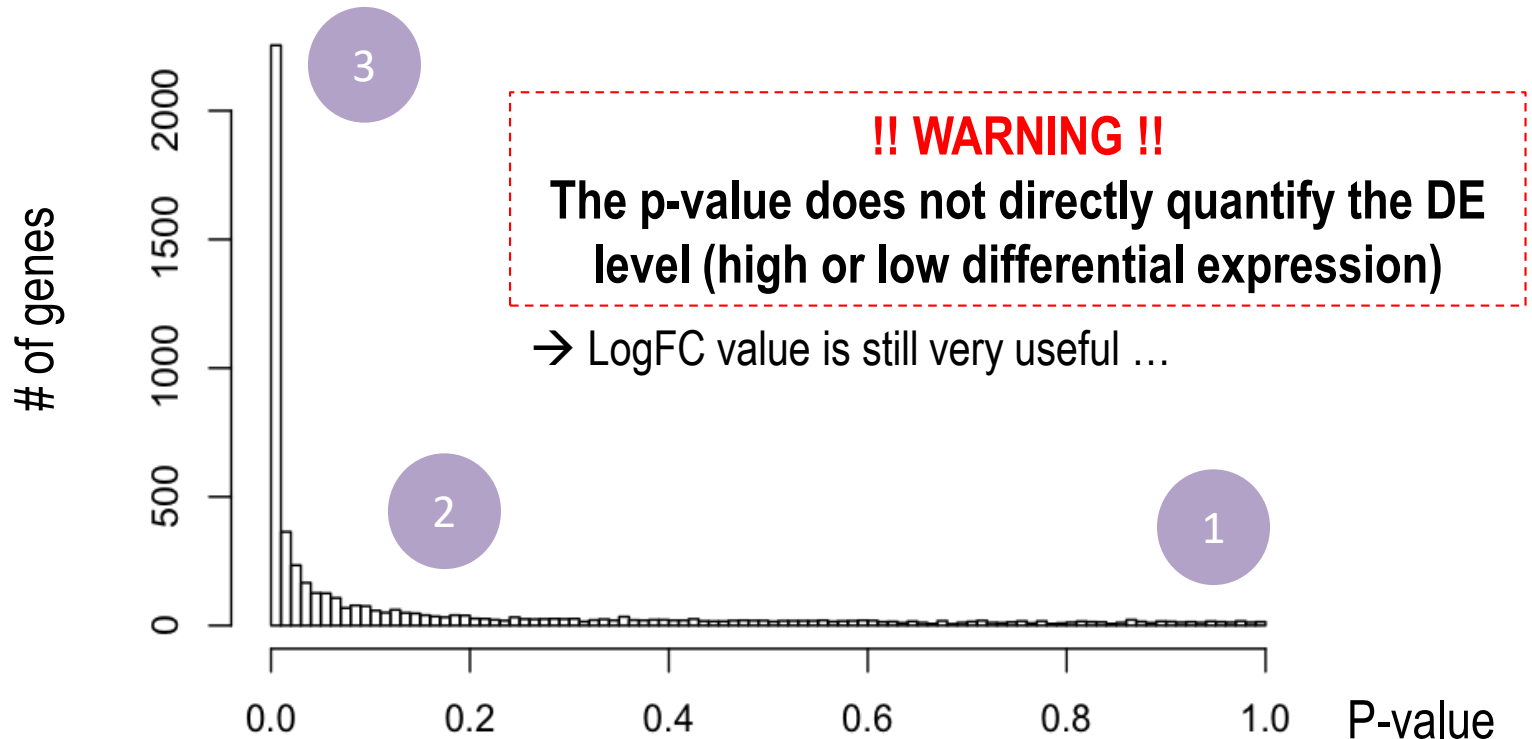
P-values to quantify confidence levels

- A p-value **quantifies the risk** to select a gene as differentially expressed, whereas it is not



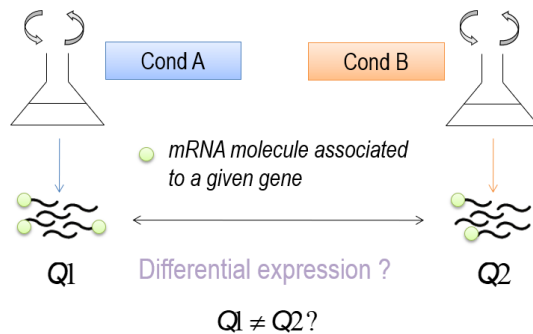
P-values to quantify confidence levels

- A p-value quantifies the risk to select a gene as differentially expressed, whereas it is not



Multiple testing correction (do you remember ?)

What is the question ?



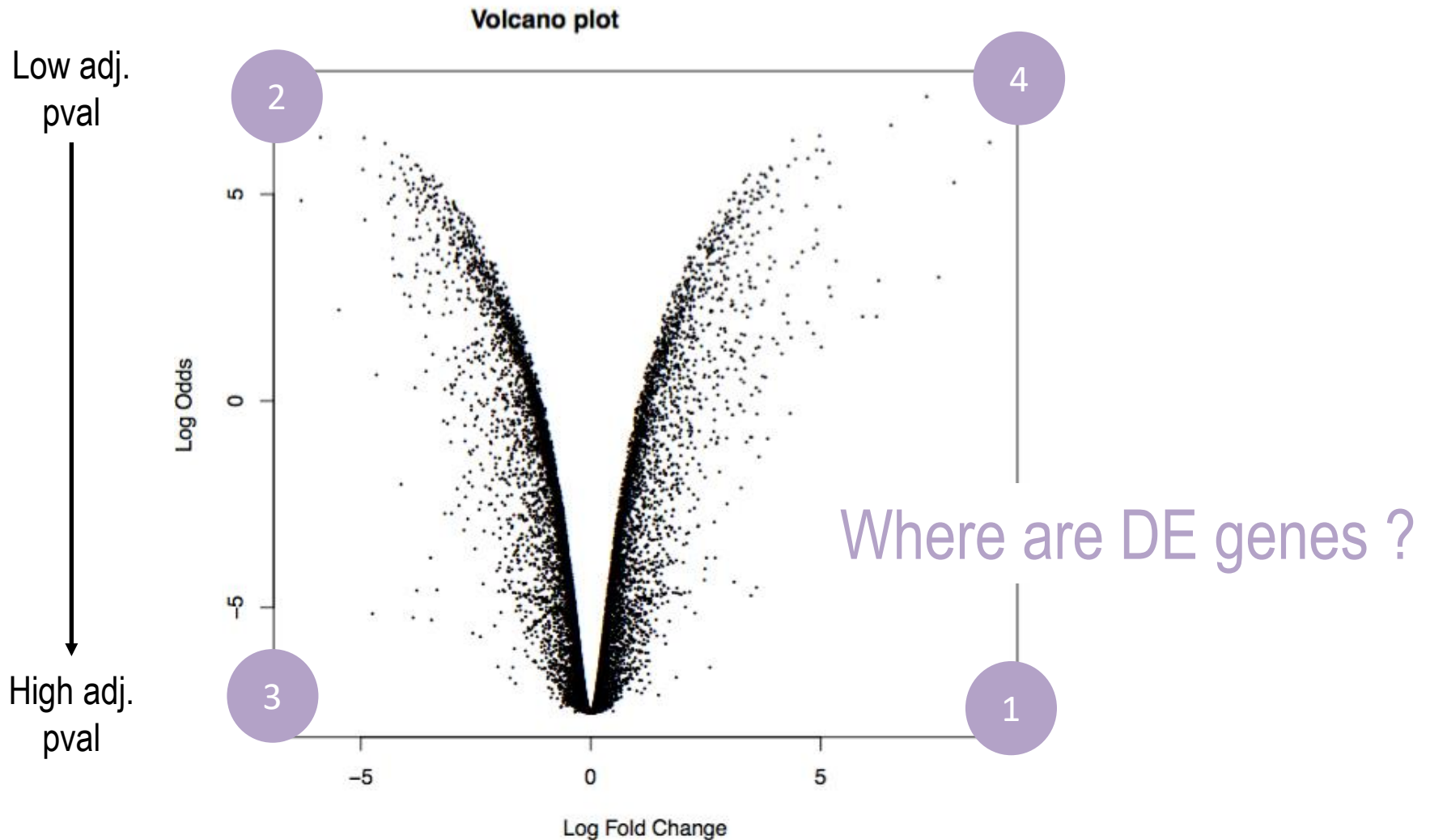
➤ All genes are successively analyzed in a differential analysis

If conditions A and B are identical (no DE gene), since many statistical tests are performed (for instance 10 000) ...

... we can expect **500 false positive** genes*, with a standard p-value cut off of 5%

* Expected value according to a binomial distribution

Combining logFC and adj. P-values ...

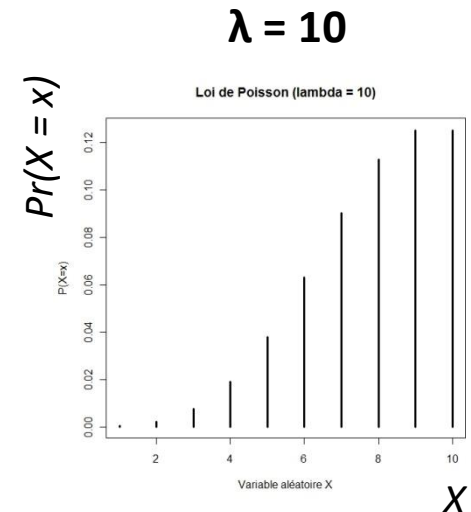
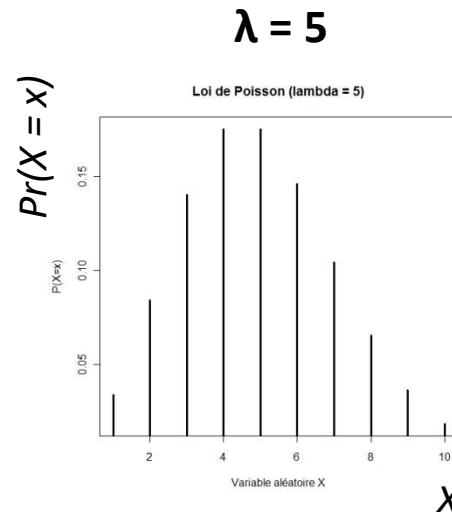
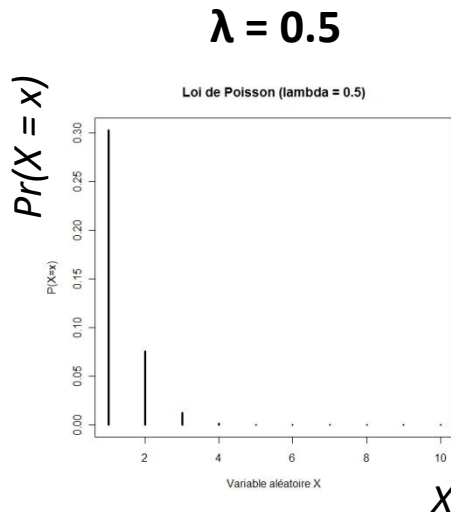


Other strategy to calculate p-values : Poisson Model

$$p(k) = P(X = k) = \frac{\lambda^k}{k!} e^{-\lambda}$$

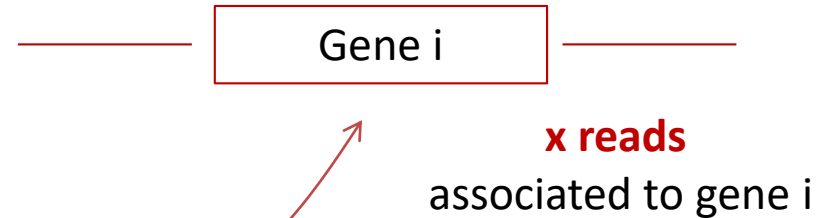
- Do you remember ?
 - Bernoulli trial
 - Binomiale distribution
 - Poisson distribution (rare events)

How to estimate the
value of λ ?



Application to HTS data analysis

Read sequences



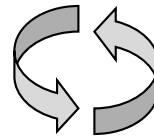
Bernoulli trial =
Sequencing of a read

p

$1 - p$

Gene i

Other gene



Repetition of Bernoulli
trials (for each read)

P and λ are unknown
parameters

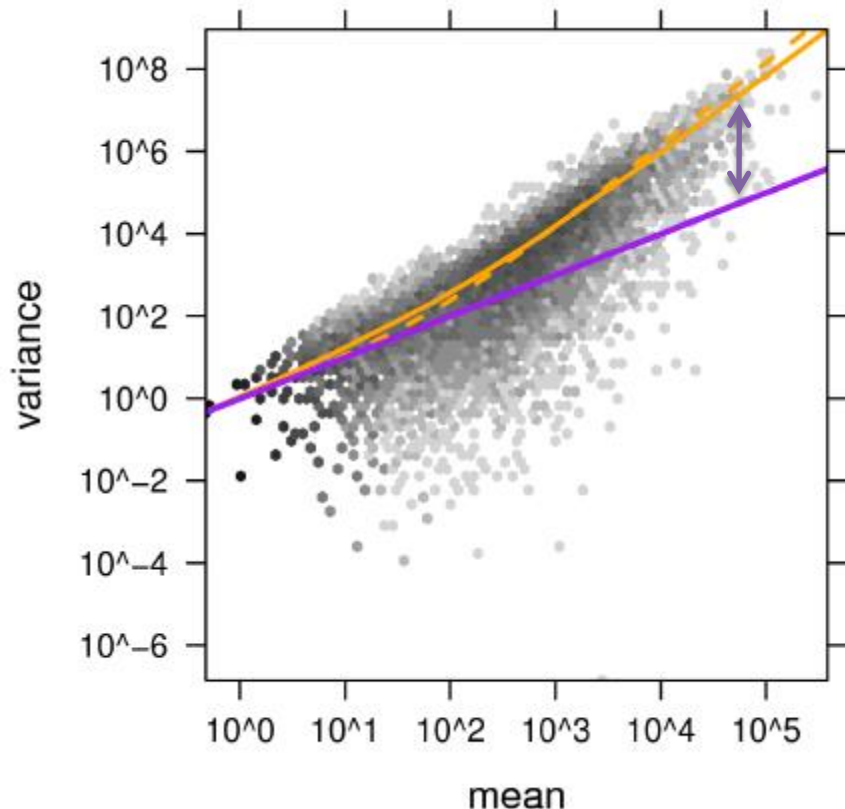
❖ The number of reads associated to the gene i is a random variable that follow a Binomial distribution $B(n, p)$.

❖ As n is high and p small, Poisson distribution can be used.

Is the Poisson Model Relevant ?

- ❖ The Poisson model has the major interest to require the estimation of only one parameter (λ).

Mean and variance calculated from two replicates



Property:
 $E(X) = Var(X) = \lambda$

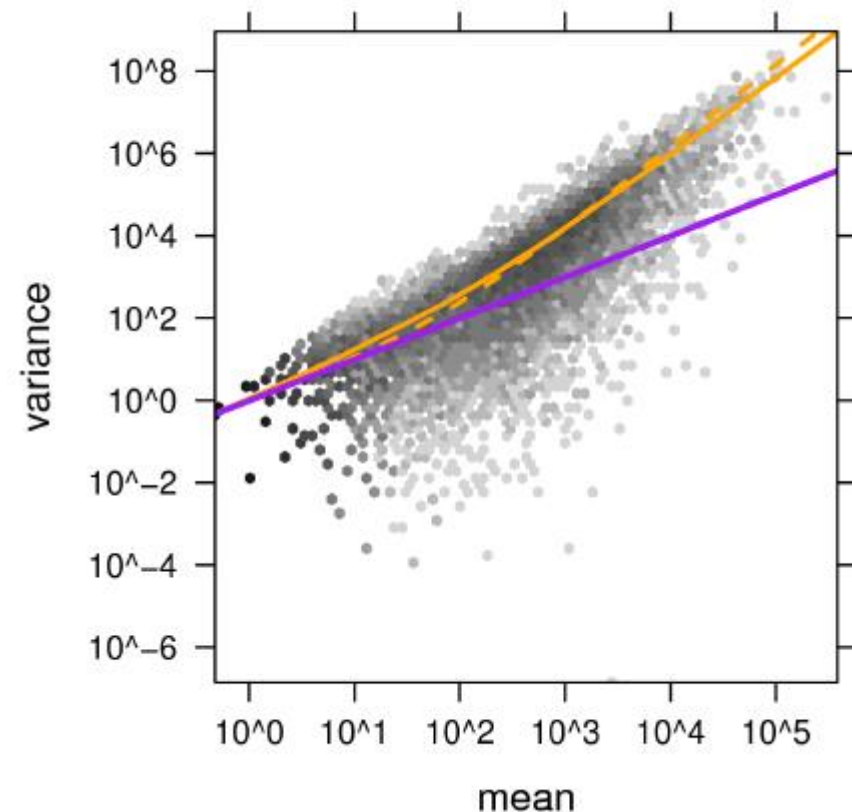
- ❖ Distribution expected if the Poisson model is true.

Variability associated to the number of reads increases more than the mean value

- ❖ Poisson model under-estimates the variability for highly expressed genes (Increases false positive rate).

Variance Estimation is “the key”

Mean and variance calculated from
two replicates



Variance calculated from comparing two replicates

Poisson

$$v = \mu$$

Poisson + constant CV

$$v = \mu + \alpha \mu^2$$

Poisson + local regression

$$v = \mu + f(\mu^2)$$

❖ **DESeq** (Anders et al. 2010):

❖ **edgeR** (Robinson et al. 2008):

Gene expression level is another important information

- Two genes can have identical logFC, but different levels of transcriptional activity

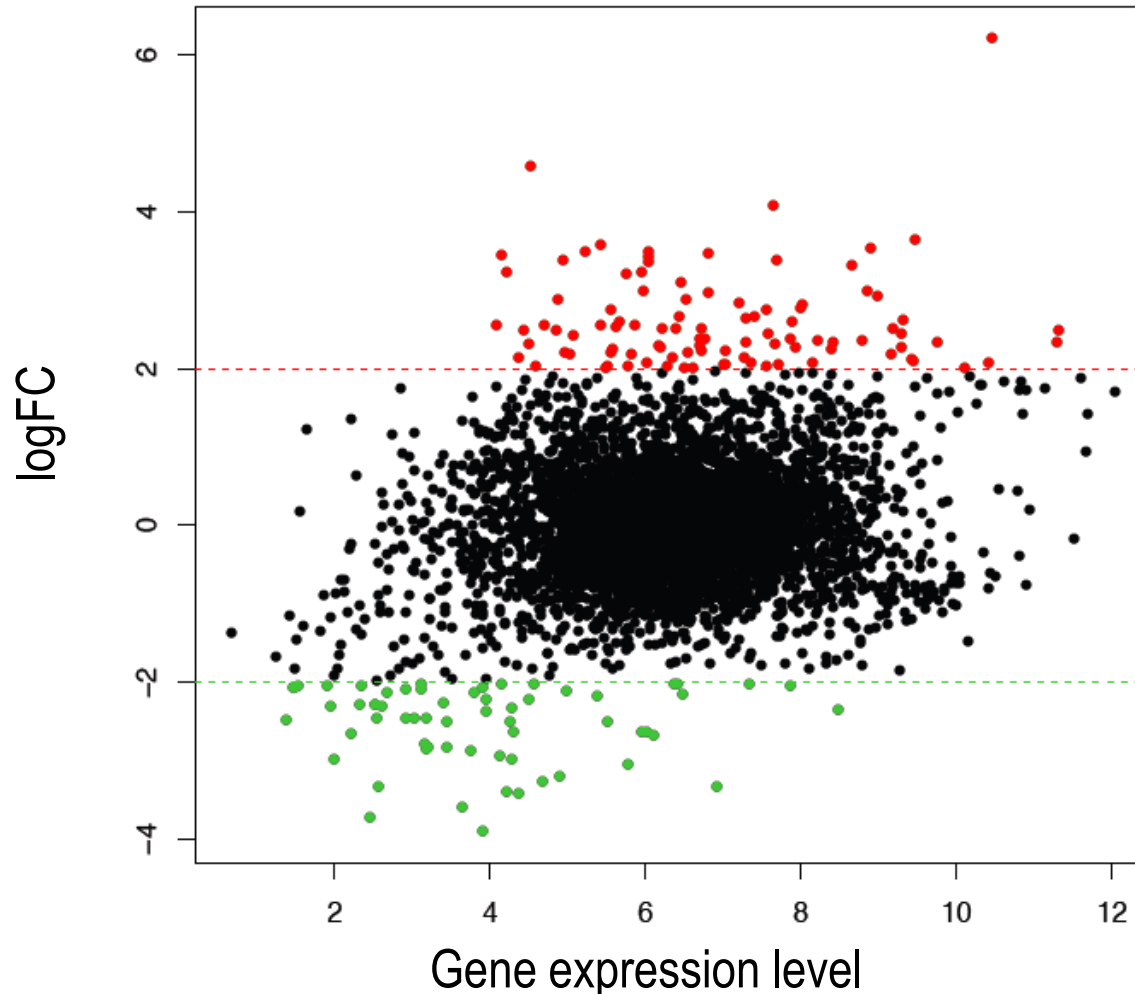
+ + + + Gene 1: $\log FC_{G1} = \log 2(8000 / 1000) = 3$

Gene 2: $\log FC_{G2} = \log 2(8 / 1) = 3$

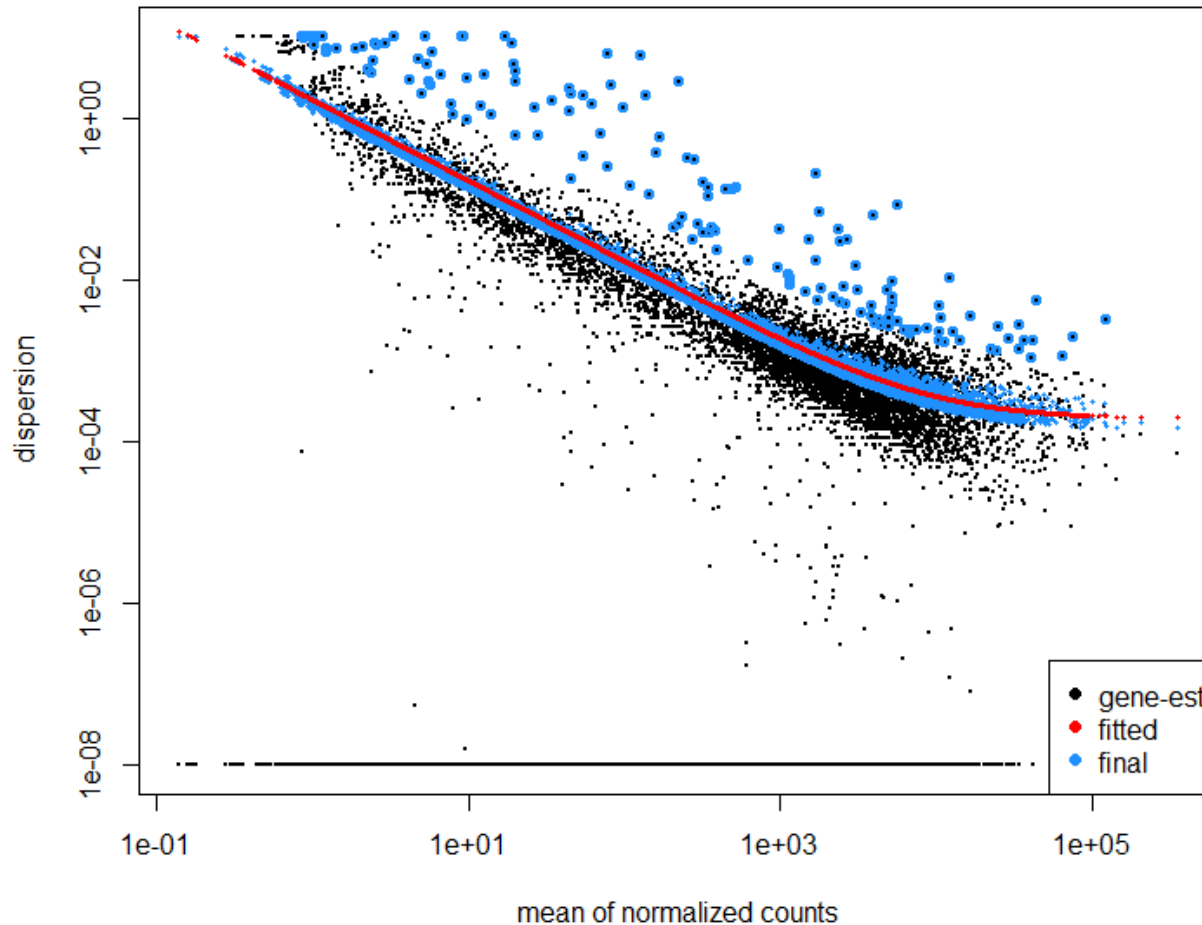
Random fluctuation ?

Same level of confidence ?

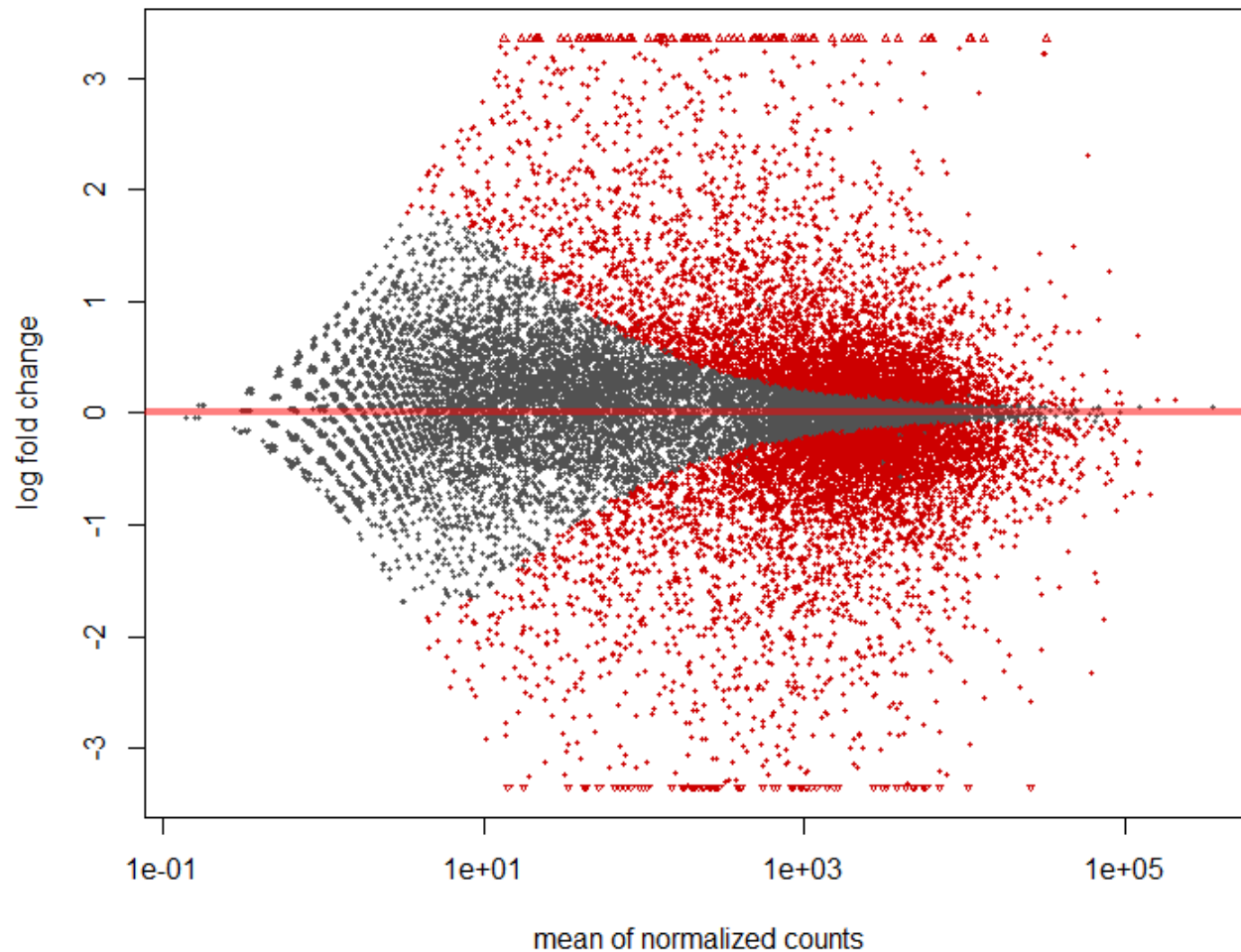
Gene expression level is another important information



Results obtained with DESeq2 (1/2)



Results obtained with DESeq2 (2/2)



R script to be used/modified

Search_DEgenes.R

```
search_DEgenes_M2Bioinfo.R
search_DEgenes_M2Bioinfo.R No Selection

#####
# Gaelle LELANDAIS <gaelle.lelandais@u-psud.fr>
# Melina GALLOPIN <melina.gallopin@u-psud.fr>
#####

####
# Ce script R est destine aux etudiants de Master. Il a pour objectif de les aider dans la
# comprehension des methodes de selection des genes differentiellement exprimes.
####

####
# Librairies R necessaires : LIMMA, MASS, DESeq2, edgeR
####

#-----
#-----

pdf("graphics.pdf")

# Lecture des données
countData = read.table("count_dataFile.txt", header = T, row.names = 1)
# --> Ces donnees sont en relation avec l'article :
#
# "Determination of a Comprehensive Alternative Splicing Regulatory Network and
# Combinatorial Regulation by Key Factors during the Epithelial-to-Mesenchymal Transition"
#
# Yang et al. (2016, Molecular and Cellular Biology)
#
# --> Deux points de temps ont ete etudies par RNAseq : "Day0" et "Day7", 3 replicats sont
# disponibles pour chaque point de temps.

#-----
# Pretraitement des donnees entre les experiences (filtrage des faibles et normalisation)
#-----

## Suppression des genes sans valeurs (0 pour toutes les experiences)
length(which(rowSums(countData)==0)) ## 2957
countData <- countData[-which(rowSums(countData) == 0),]

# Distribution des valeurs de comptages dans les différents échantillons
boxplot(log(countData + 1), ylab = "#reads (log scale)",
        main = "Boxplot of read counts", col = c(rep("grey", 3), rep("blue", 3)),
        names = c(paste("Day0_", 1:3, sep = ""), paste("Day0_", 1:3, sep = "")))
# --> Les distributions des donnees de comptage sont tres proches entre les differentes
# conditions.
```

To conclude,

Statistics help you to take a decision (is the gene differentially expressed ?) ...

... it should not take the decision for you

➤ **Important parameters to consider in a DE analysis:**

Differential expression

Confidence level

Gene information in
databases, literature ...

Data reproducibility

Gene expression level

SECTION 3

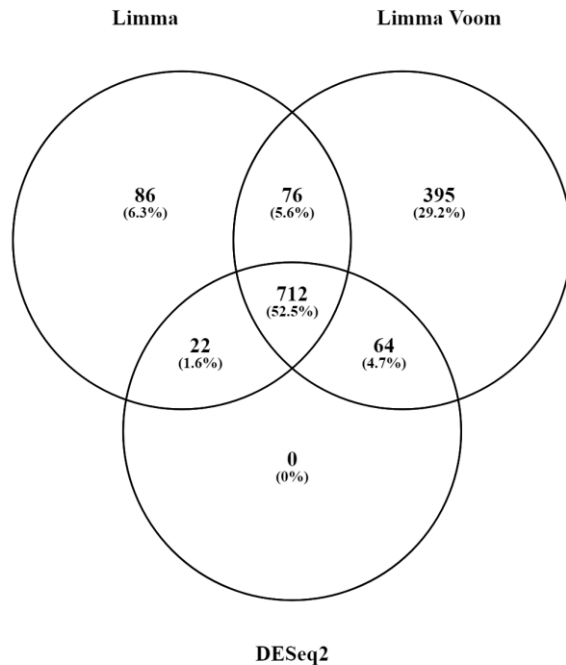
PRACTICAL TRAINING

Your work ...

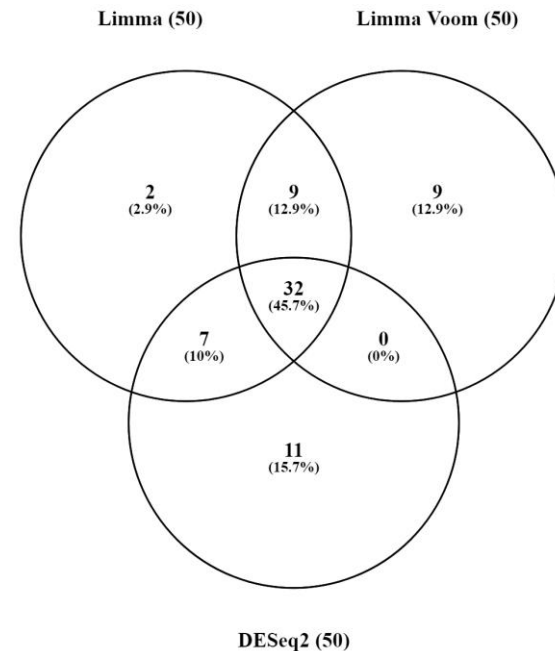
- Compare results obtained with different approaches (logFC threshold, DESeq2, etc.)
- Search for genes located on Chr18 in your lists of DE genes. Check for read coverage with IGV
- Define (or refine) your own list of “favorite genes”

Method comparisons

- Genes with logFC values > 2 or < -2 AND adj.pal < 0.01 were selected and compared between the three methods :

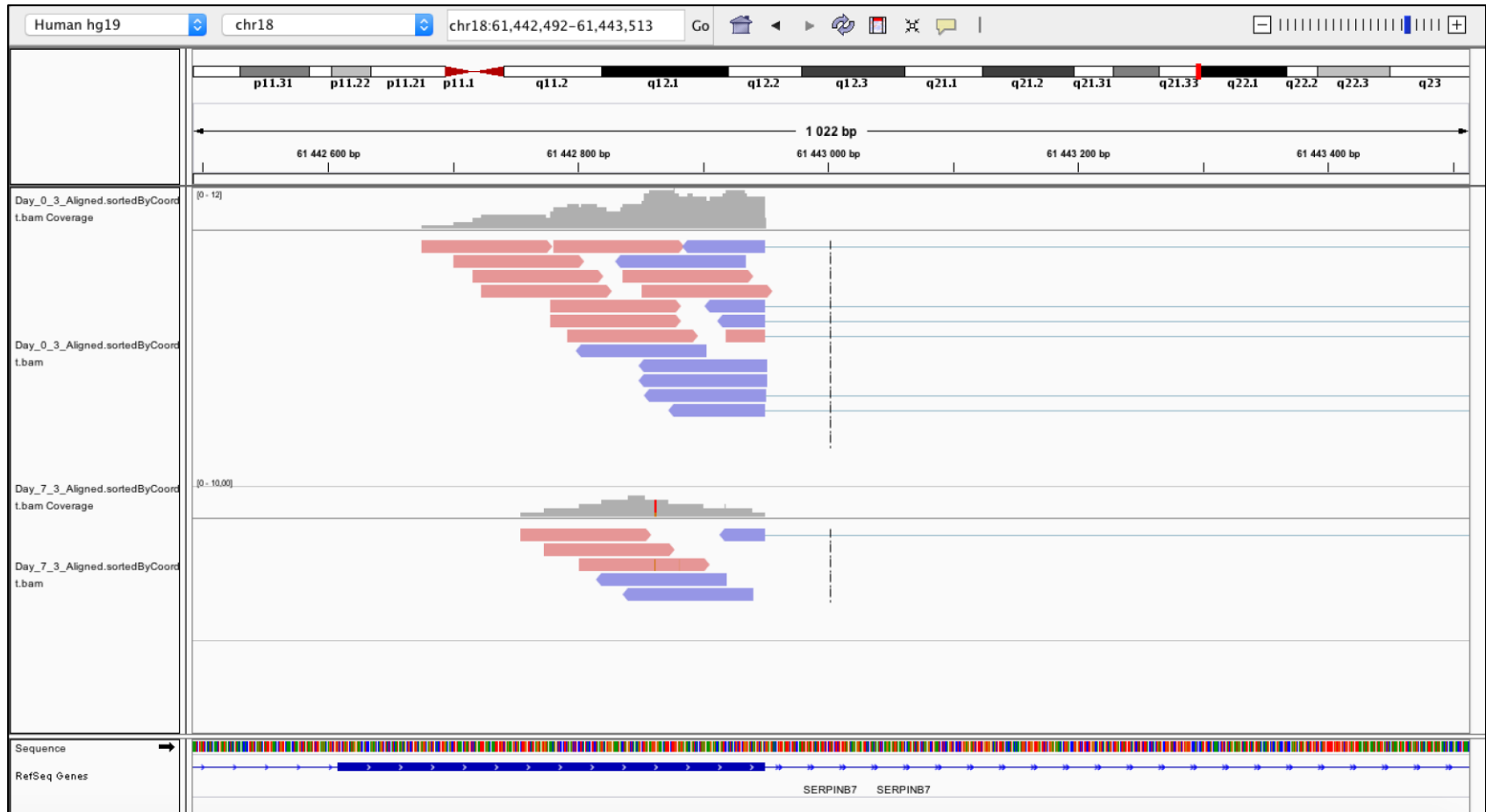


All selected genes



Top 50 genes

DE genes verified with IGV (1/2)



DE genes verified with IGV (2/2)

