



Session 3

Statistiques pour les données omiques

Teachers: Claire Vandiedonck, Jacques van Helden

Helpers: Antoine Bridier-Nahmias, Anne Badel

Plan de la séance

Retour sur les séances 1 et 2:

- debrief sur les commandes R
- TP - part I : données simulées
- debrief sur les stats de base

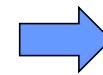
Coffee break

Statistiques pour les données omiques:

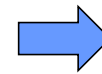
- TP – part II : "industrialisation" des tests d'hypothèses
- cours – part I :
 - donner du sens aux données omiques et problèmes de dimensionnalité
 - 1^{er} problème: tests multiples
- TP –part III: tests multiples
- cours – part II :
 - 2^{ème} problème: estimation des paramètres des distributions
 - 3^{ème} problème: réduction de la dimensionnalité -> cf. sessions suivantes

Liens

Deux difficultés dans la mise en evidence d'effet



grande masse
de données



issues d'échantillons
et non de la population
en partie cachée

1. Introduction: making sense of omic's data

COMMENTARY

'Ome Sweet 'Omics-- A Genealogical Treasury of Words

By Joshua Lederberg and Alexa T. McCray

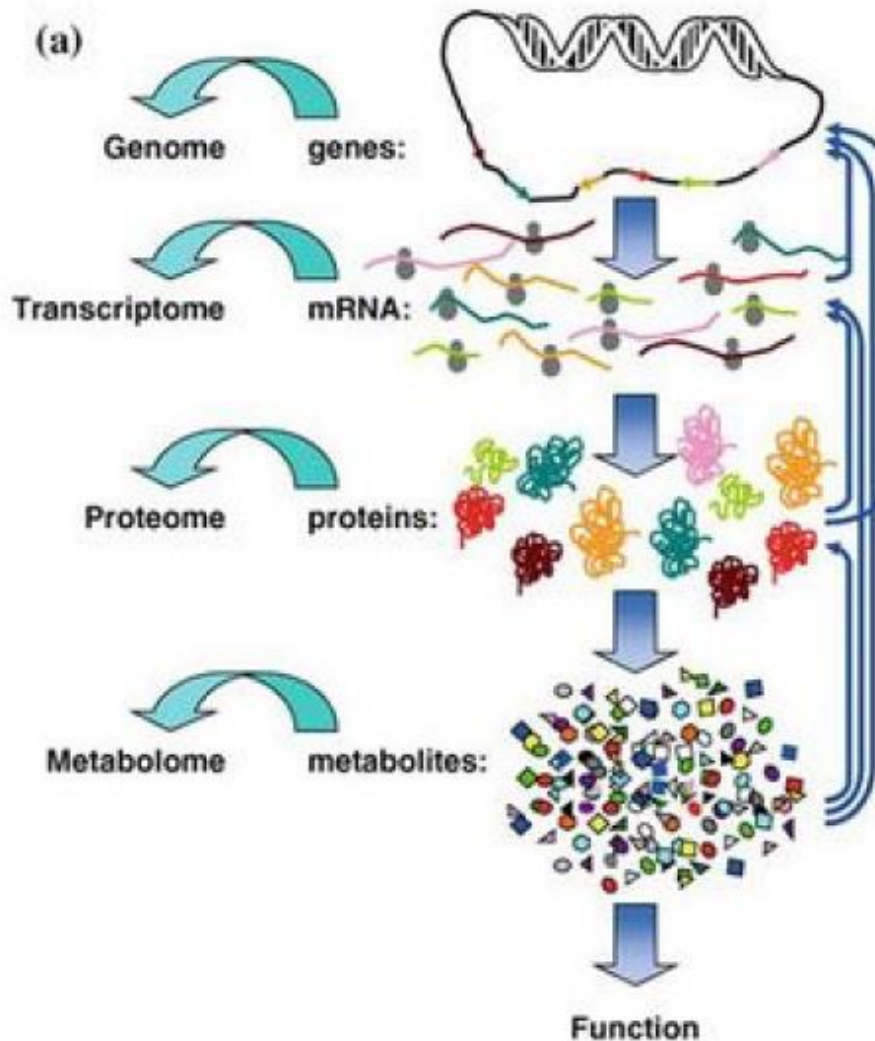
The Scientist 15[7]:8, Apr. 2, 2001

| | | |
|------------|----------------|---------------|
| antigenome | immunogenome | plastidome |
| bacteriome | immunome | plerome |
| basidiome | haptenome | proteinome |
| biome | karyome | proteome |
| cardiome | leptome | psychome |
| caulome | microbiome | regulome |
| chondriome | mnemome | rhabdome |
| cladome | mycetome | rhizome |
| coelome | neurome | stereome |
| epigenome | odontome | thallome |
| erythrome | osteome | tracheome |
| genome | pharmacogenome | transcriptome |
| geome | phenome | trichome |
| hadrome | phyllome | vacuome |
| histome | physiome | |

Genomics and *Proteomics* are the buzzwords of the dawning millennium. There is no counting of www.ics.com and www.ix.com sites to be found on the Web. That most of these terms, old and new, have been contrived as slogans to attract attention, does not diminish their likely substance, and they are embedded in the advancing edge of science and technology.

<https://lhncbc.nlm.nih.gov/system/files/pub2001047.pdf>

Integration des données omiques



The “Omics” Cascade

What can happen

GENOME



What appears to be happening

TRANSCRIPTOME



What makes it happen

PROTEOME



What has happened and is happening

METABOLOME



PHENOTYPE

Hétérogénéité des données omiques

Nature des données

- binaires (eg. présence ou absence d'un allèle ou d'un site de liaison)
- catégoriques (séquences de site consensus, isoforme exprimée)
- quantitative discrète (génotypes: 0, 1, 2)
- quantitative continue (niveau d'expression d'un gène ou d'une protéine)

Dimension des données (*exemples chez l'homme*)

- génome (4×10^6 de variants bi-alléliques de type SNP)
- transcriptome (20-60 000 gènes, 200 000 transcrits)
- protéome (18 000 protéines, 293 000 peptides)

Données manquantes (4000 protéines)

Structure des données

- corrélations entre les variables mesurées (déséquilibre de liaison, co-expression...)
- corrélations entre les types de données

En plus, des données non-omiques peuvent exister = co-variables

| | G_1 | G_2 | ... | G_p | condition | age | gender | BMI | glycemia |
|---------|-------|-------|-----|-------|-----------|-----|--------|-----|----------|
| $i = 1$ | 0 | 12 | | 41 | healthy | 38 | W | 22 | 0.8 |
| $i = 2$ | 10 | 3 | | 2 | affected | 15 | M | 30 | 0.2 |
| . | | | | | | | | | |
| . | | | | | | | | | |
| $i = N$ | 0 | 20 | | 15 | affected | 90 | W | 31 | 1.5 |

omics data

facteur d'intérêt qu'on veut tester

covariables (metadata)

samples

- Par exemple, on peut avoir le niveau d'expression par gène pour chaque échantillon
 - On peut aussi avoir des données cliniques pour les échantillons incluant le facteur d'intérêt qu'on veut tester et d'autres covariables qui pourraient impacter les niveaux d'expression
- On souhaite expliquer les variations d'expression (variable expliquée) en fonction de covariables cliniques (variables explicatives)

Why using statistics ?

Making sense of data

- ↳ **Aim:** identify variables whose variation levels are associated with a phenotype or a covariate of interest
(eg: response to stress, to a treatment, survival, mutation, tumor class, time...)

Variable to explain \sim explanatory variables + covariates + residual error

Problems addressed by statistics:

1. **estimation:** of the effects of interest and of how they vary
2. **testing:** = assessing the statistical significance of the observed effects

Quels facteurs peuvent expliquer la variation d'un trait?

Variation inter-groupes

1. Facteur/covariables d'intérêt => design experimental

- ✓ conditions expérimentales testées: stimulus, traitement, temps, maladie...
- ✓ variabilité génétique: mutation
- ✓ tissus/type cellulaire...

2. Variation technique: répliqués techniques

- ✓ experimental: lot, jour, expérimentateur, température ambiante...
- ✓ multiplexage
- ✓ variation de plate-forme

Variation intra-groupes

Variation biologique => répliqués biologiques

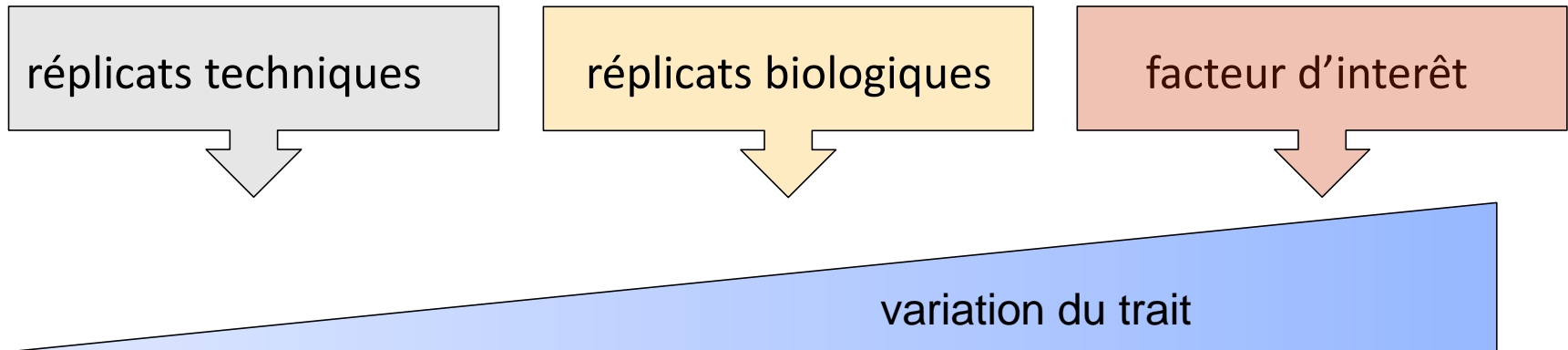
- ✓ fluctuation d'échantillonnage

De l'importance d'un bon design experimental

Les différences entre les conditions peuvent uniquement être testées si des **REPLICATS** sont inclus

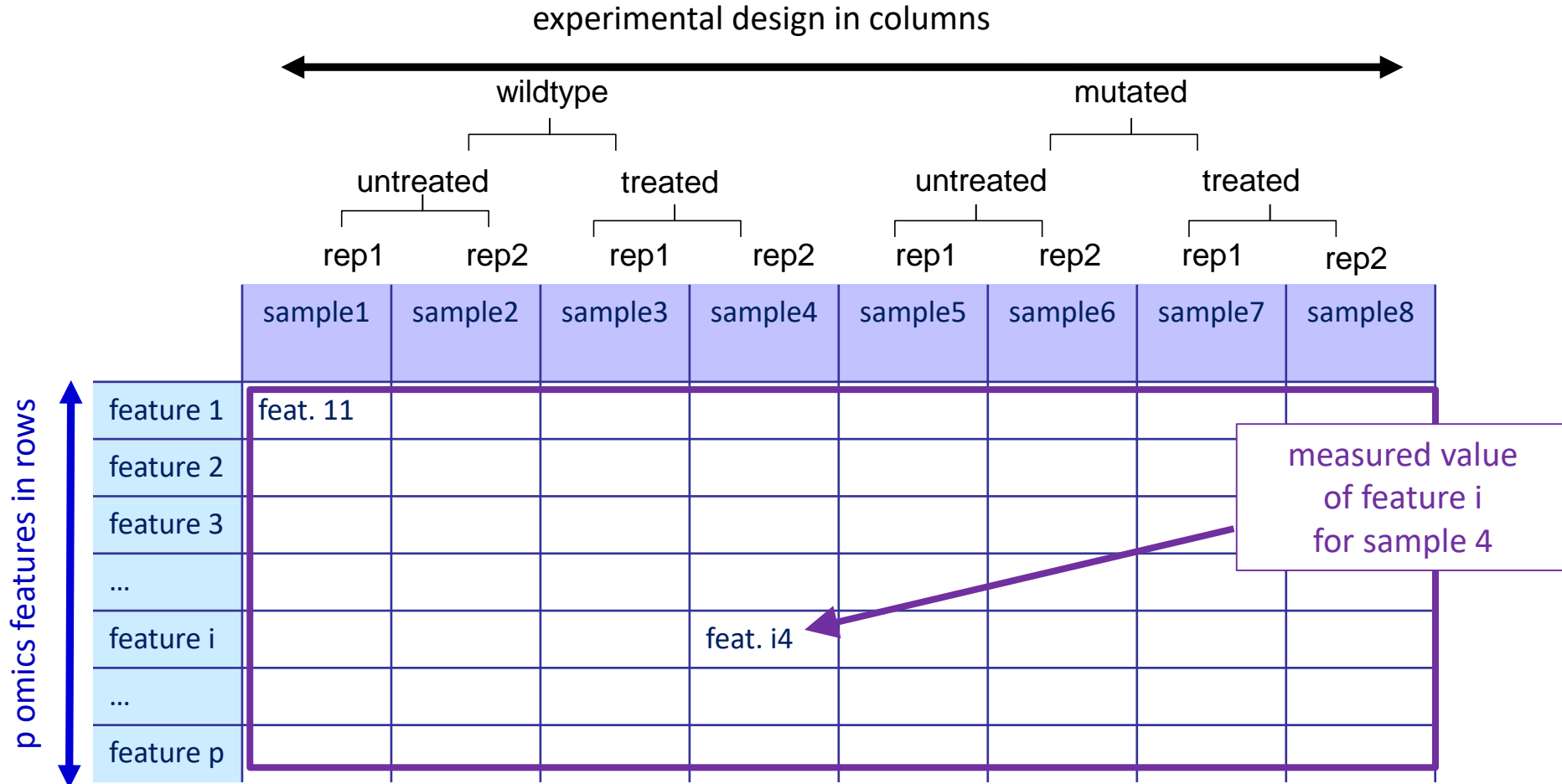
⇒ permettent de déterminer quelles différences sont dues aux fluctuations aléatoires d'échantillonnage

👉 **Ideal scenario :**



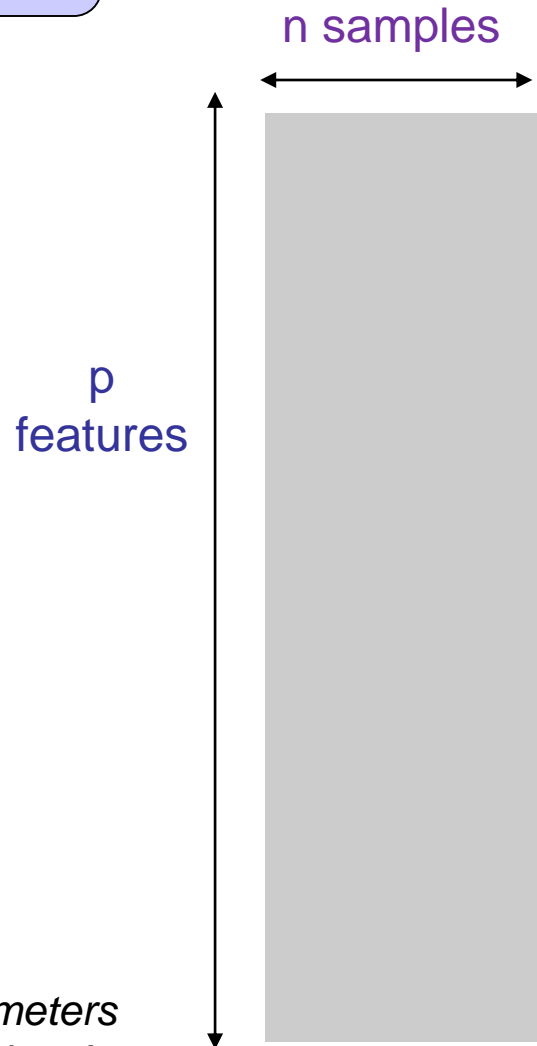
La structure des données omiques

Matrice de données



Les problèmes de dimensionnalité

$$p \gg n$$



p = number of parameters (features), not p-values!

Les problèmes de dimensionnalité

$$p \gg n$$

↪ 3 problèmes

p
features

n samples



n small:

↪ difficulty to estimate parameters of each trait distribution

p = number of parameters (features), not p-values!

Les problèmes de dimensionnalité

$p \gg n$

↪ 3 problèmes

p
features

p large:

↪ multiple testing issue

n samples

n small:

↪ difficulty to estimate parameters of each trait distribution

p = number of parameters (features), not p -values!

Les problèmes de dimensionnalité

$p \gg n$

↪ 3 problèmes

p
features

p large:

↪ multiple testing issue

n samples

n small:

↪ difficulty to estimate parameters of each trait distribution

Correlation between traits

↪ difficulty to estimate because n small
↪ redundancy: too many tests?

p = number of parameters (features), not p -values!

2. The 1st issue: multiple testing

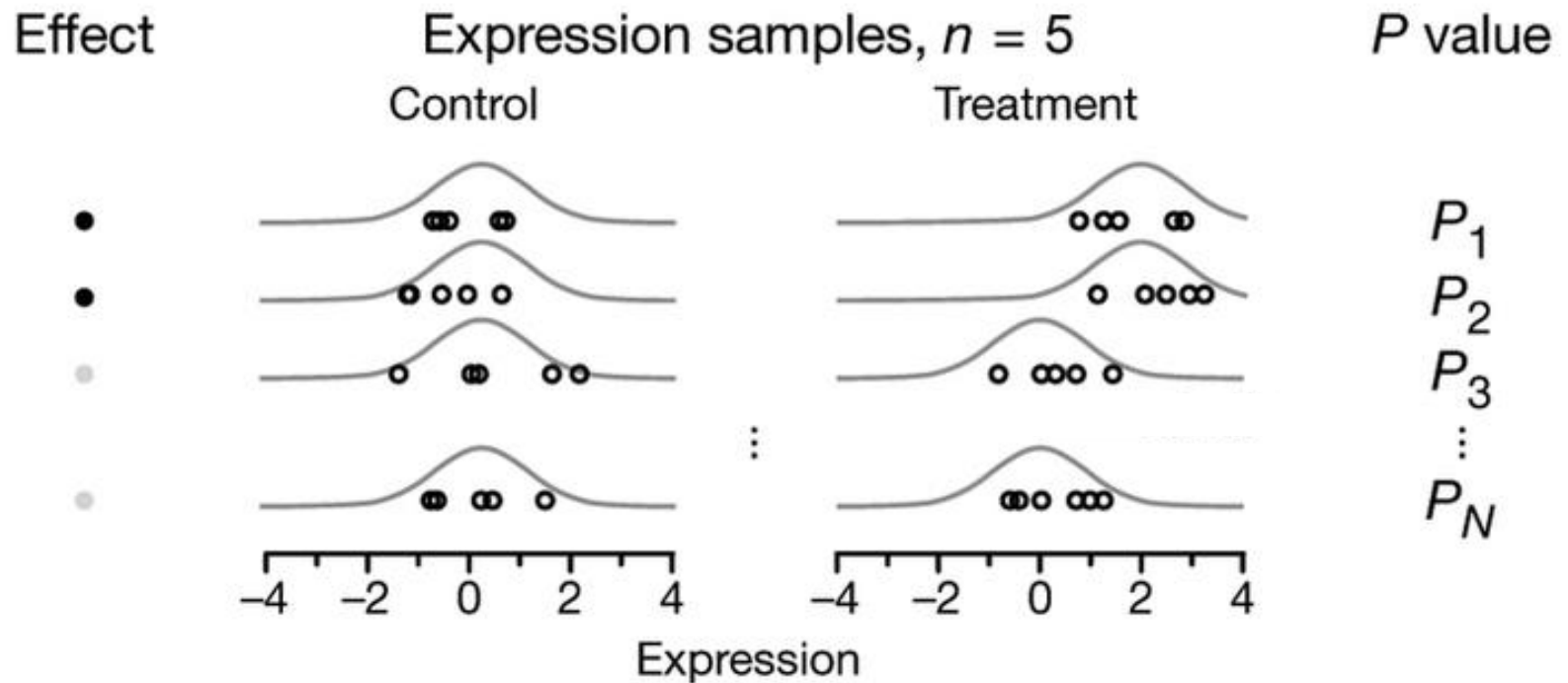
The problem

We perform multiple tests = one per feature/trait

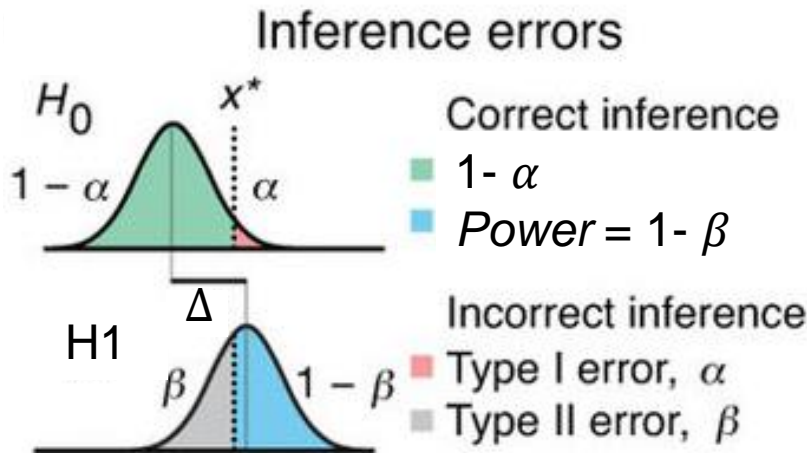
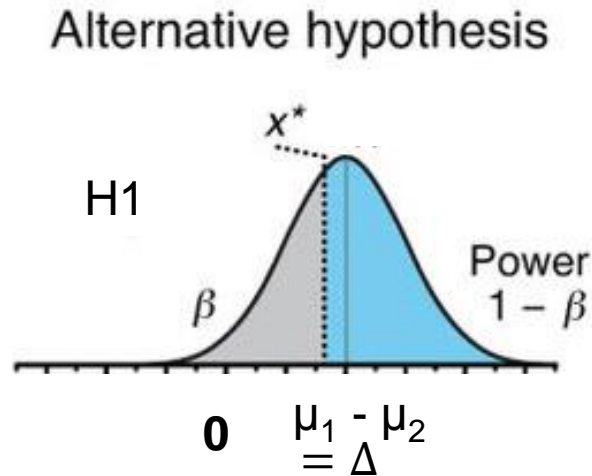
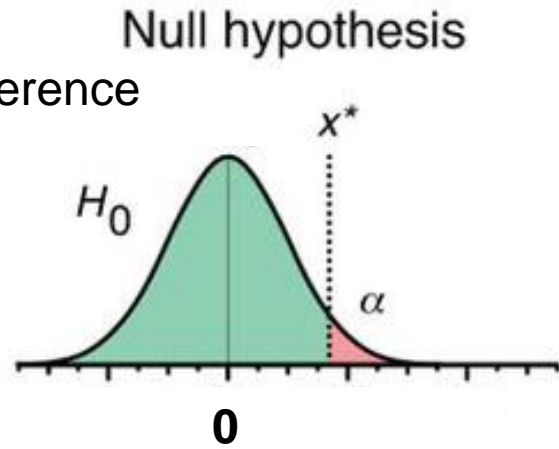
↳ for each feature, we either reject or not H_0 at a risk $\alpha = \text{PCER}$

= per-comparison error rate

Simulation gene expression samples



Test theory : alpha and beta risks



| Test decision | Reality | |
|--------------------|----------------------|---------------------|
| | H_0 | H_1 |
| no reject of H_0 | $1 - \alpha$ (TN) | β (FN) |
| reject of H_0 | α (FP) | $1 - \beta$ (TP) |

Why is the problem so important?

Omics are big data:

A typical microarray or RNA-seq experiment: 10,000 genes

=> as many hypothesis tests

Just one hypothesis test:

For an $\alpha = 0.05$, we tolerate to reject H_0 wrongly 5% of the times

↳ but for 10,000 tests the number of false positives goes up to 500

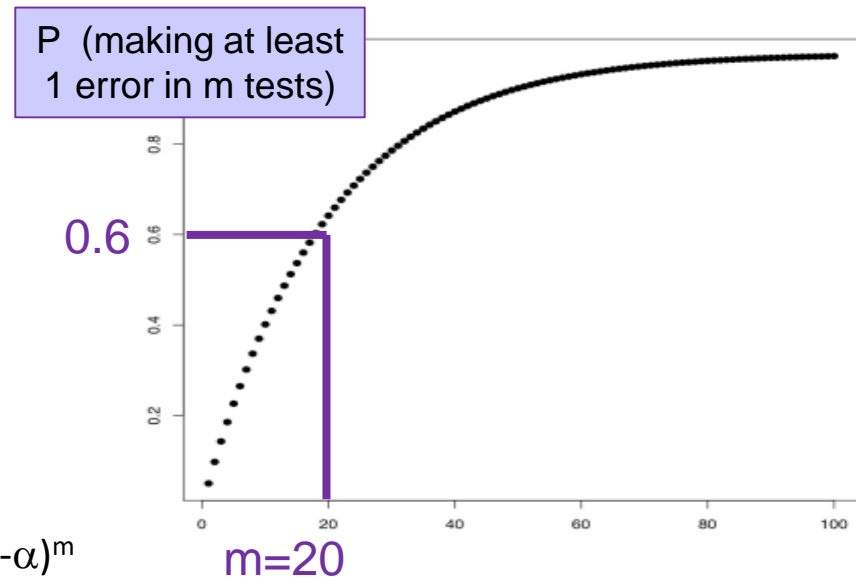
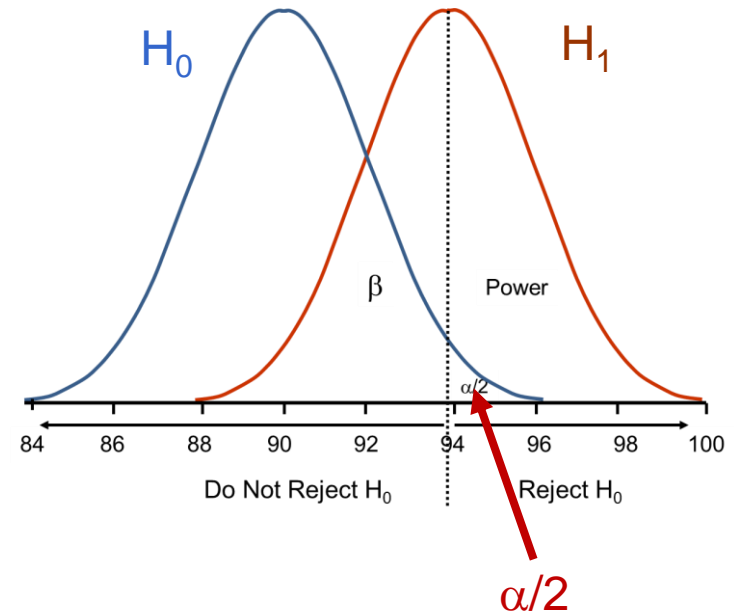
=> too many!!!

Expected value (**e-value**)

- Expected number of FP = $E(\text{FP}) = m\alpha$

Family-wise error rate (FWER)

- $P(\text{making an error}) = \alpha$
- $P(\text{not making an error}) = 1 - \alpha$
- $P(\text{not making an error in } m \text{ tests}) = (1 - \alpha)^m$
- **FWER** = $P(\text{making at least 1 error in } m \text{ tests}) = 1 - (1 - \alpha)^m$



Counting errors

| Decision on H_0 | H_0 True | H_1 True | |
|-------------------|-------------------------|-------------------------|-----|
| reject | V (incorrect) | S | R |
| do not reject | U | T (incorrect) | m-R |
| | m_0 | $m - m_0$ | m |

m = number of tests

R = number of rejected H_0

m_0 = number of true H_0

➤ only m and R are observed!

V = number of type I errors = **false positives**

By the way, where are:

the false negatives?

the **true positives**?

the **true negatives**?

Counting errors

| Decision on H_0 | H_0 True | H_1 True | |
|-------------------|-------------------------|-------------------------|-----|
| reject | V (incorrect) | S | R |
| do not reject | U | T (incorrect) | m-R |
| | m_0 | $m - m_0$ | m |

m = number of tests

R = number of rejected H_0

m_0 = number of true H_0

➤ only m and R are observed!

V = number of type I errors = **false positives**

By the way, where are:

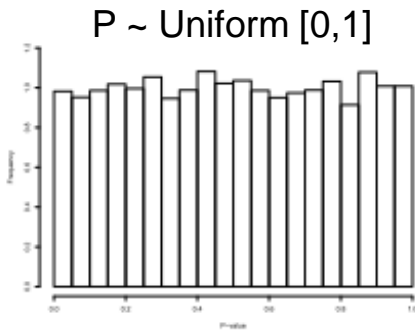
the false negatives?

the **true positives**?

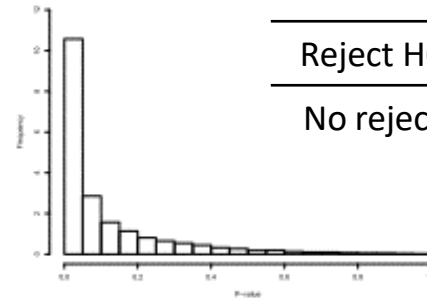
the **true negatives**?

| | H_0 True | H_1 True |
|--------------|------------|------------|
| Reject H_0 | FP | TP |
| No reject | TN | FN |

Controlling the type I error rate



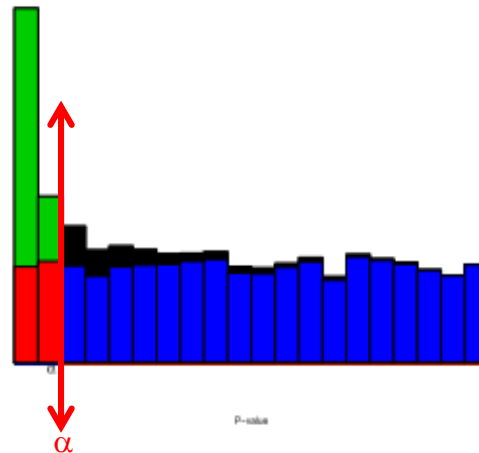
p values for null genes



p values for non-null genes

| | H0 True | H1 True |
|-----------|---------|---------|
| Reject H0 | FP | TP |
| No reject | TN | FN |

p values for all genes



True positives = S
 False positives = V
 True negatives = U
 False negatives = T

Where to set the threshold of significance to control the type I error rate?

=> Trade-off between type I error and power!!

Bonferroni correction

Aim: to control the family-wise error rate (FWER):

= the error rate across the whole collection/family of hypothesis tests

= FWER = $P(V \geq 1)$ = probability of ≥ 1 false positive among all tests

↪ By “adjusting” the p value with the Bonferroni correction

set $\alpha' = \alpha/m$
reject hypotheses if $p < \alpha'$

✓ E.g. for a type I error rate of 0.05 per experiment (PCER)
and $m = 10\,000$ tests: $\alpha' = 0.05/10,000 = 5 \times 10^{-6}$

very popular

the problem for “Omics” experiments: very conservative

=> alternative approaches investigated: very active area of current research in statistics!

False discovery rate (FDR)

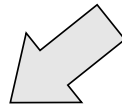
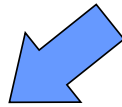
We focus on positive tests (H_0 rejected):

FDR = proportion of false positive among the set of rejected hypotheses (the “discoveries”):

✓ $FDR = V/R$

A related parameter = the False Positive Rate (FPR)

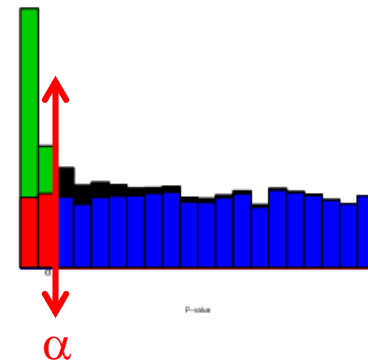
✓ $FPR = V/m_0$



| Decision on H_0 | H_0 True | H_1 True | |
|-------------------|----------------------|----------------------|-----|
| reject | V (incorrect) | S | R |
| do not reject | U | T (incorrect) | m-R |
| | m_0 | $m-m_0$ | m |

| Decision on H_0 | H_0 True | H_1 True | |
|-------------------|----------------------|----------------------|-----|
| reject | V (incorrect) | S | R |
| do not reject | U | T (incorrect) | m-R |
| | m_0 | $m-m_0$ | m |

| | H_0 True | H_1 True |
|--------------|------------|------------|
| Reject H_0 | FP | TP |
| No reject | TN | FN |



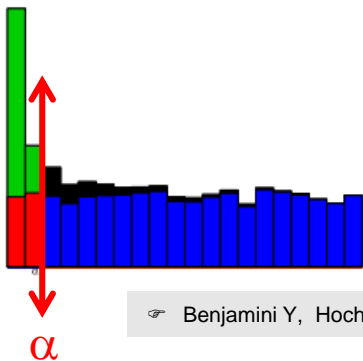
Benjamini-Hochberg procedure

To control FDR at level δ :

- ✓ order the unadjusted p-values: $p_1 < p_2 < \dots < p_m$
- ✓ find the test with the highest rank, j , for which the p value,

$$p_j \leq \delta \frac{j}{m}$$

- ✓ Declare the tests of rank $\leq j$ as significant



Example: $m = 10$
and $\delta = 0.05$

| Rank (j) | P-value | $(j/m) \times \delta$ | Reject H_0 ? |
|----------|---------|-----------------------|----------------|
| 1 | 0.0008 | 0.005 | 1 |
| 2 | 0.009 | 0.010 | 1 |
| 3 | 0.018 | 0.015 | 0 |
| 4 | 0.030 | 0.020 | 0 |
| 5 | 0.032 | 0.025 | 0 |
| 6 | 0.048 | 0.030 | 0 |
| 7 | 0.350 | 0.035 | 0 |
| 8 | 0.781 | 0.040 | 0 |
| 9 | 0.900 | 0.045 | 0 |
| 10 | 0.993 | 0.050 | 0 |

Adj. P val



$p_j \times m / j$

0.008

0.045

0.06

0.075

0.064

0.08

0.5

0.976

1

0.993

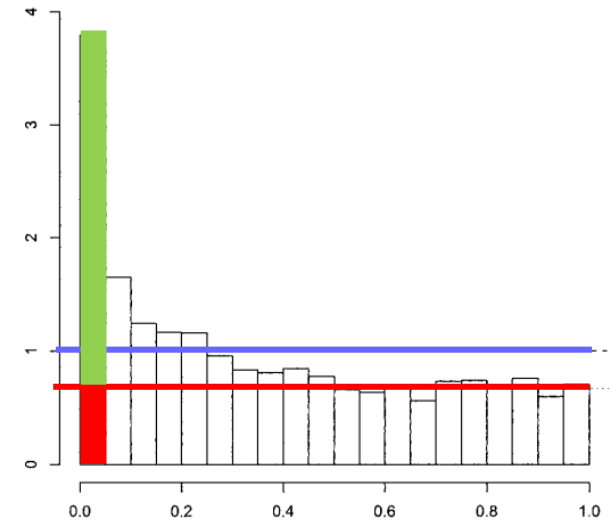
Values expected for a uniform distribution of p_j between 0 and delta

Q values

Qvalue of a gene = expected proportion of false positives when calling that gene significant

- ✓ the q-value depends on the p-value for the test of the gene and on the distribution of the entire set of p-values from the family of tests being considered (Storey and Tibshiriani 2003)
- ✓ Thus, in a microarray study testing for differential expression, if gene X has a q-value of 0.013 it means that 1.3% of genes that show p-values at least as small as gene X are false positives
- ✓ The maths:
 - π_0 : the proportion of true null tests
 - $\alpha m \pi_0$: the number of false positives
 - $\alpha m \pi_0 / R$: an estimate of the FDR

| | H0 True | H1 True | |
|-----------|---------|-----------|-----|
| Reject H0 | FP | TP | R |
| No reject | TN | FN | m-R |
| | m_0 | $m - m_0$ | m |



— histogram expected if all genes were "null", not differentially expressed

— estimate of the proportion of true "null" p-values = π_0

■ false positives ■ true positives

3. The 2nd issue: estimation of traits distribution (mean and variance)

To estimate or not to estimate?

1. No estimation when using non-parametric tests

- less power if data fit with parametric distribution
- not suitable for designs with several factors

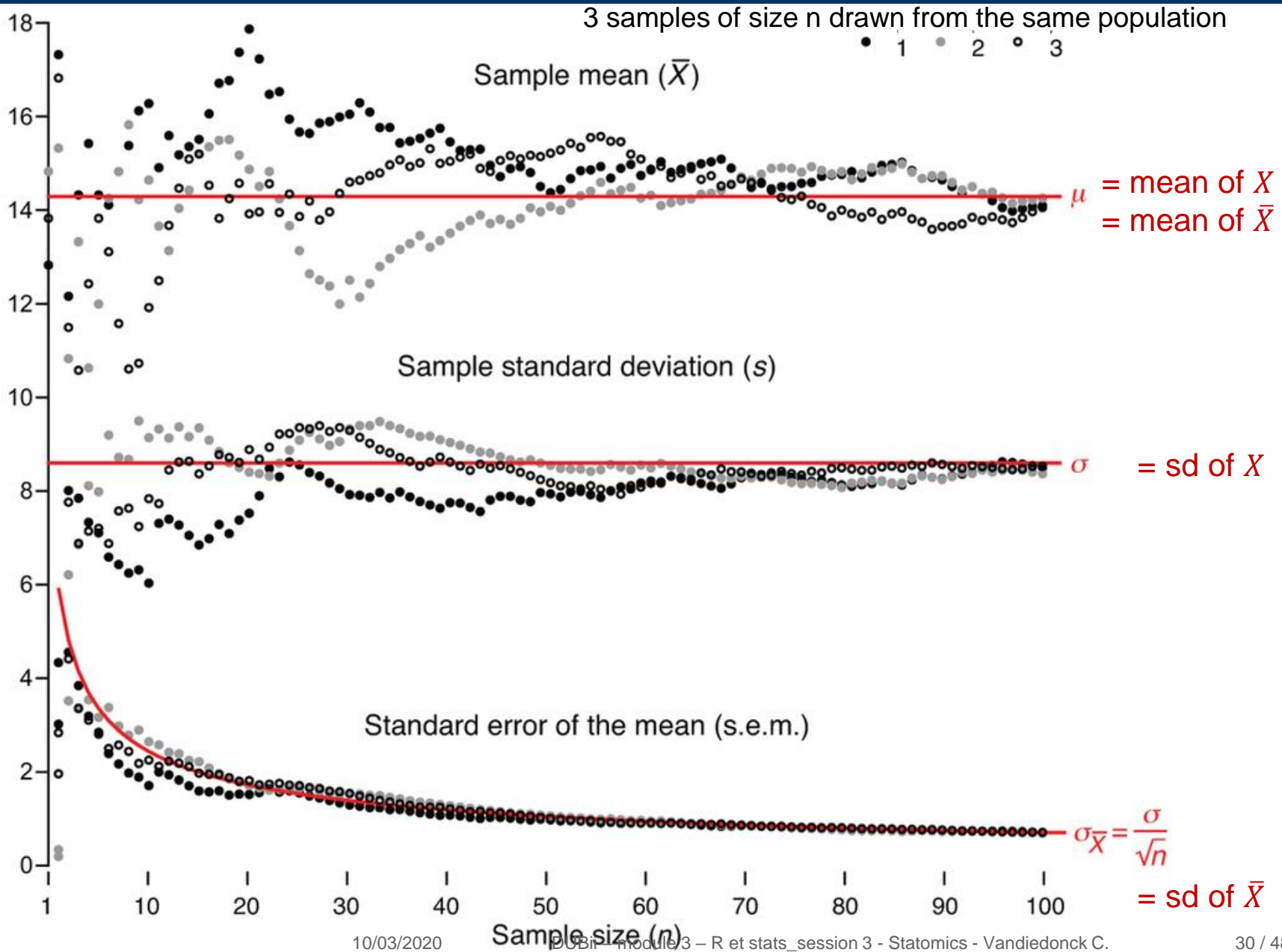
2. Random re-sampling

- approaching the distribution of p-values/statistics under null hypothesis by **permutation** (no replacement) of the levels of the factor of interest in the dataset => the empirical pvalue is the probability of observing the pvalue/statistic under the empirical distribution (cannot be lower than 1/1000 if 1000 permutations)
- estimating the CI of the distribution parameters by **bootstrap** (replacement) of the quantitated trait among all observed values within the dataset without changing the levels of the factor of interest
 - computationally intensive

3. Selecting a distribution law fitting the data

- estimation of mean and variance
- parametric tests

Better estimation when sample size is increased



Nature des données d'expression du transcriptome

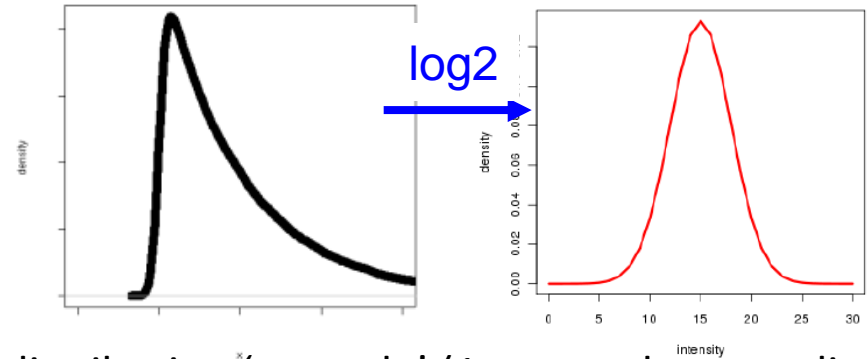
Puces

L'abondance de chaque transcript dépend de l'intensité de fluorescence

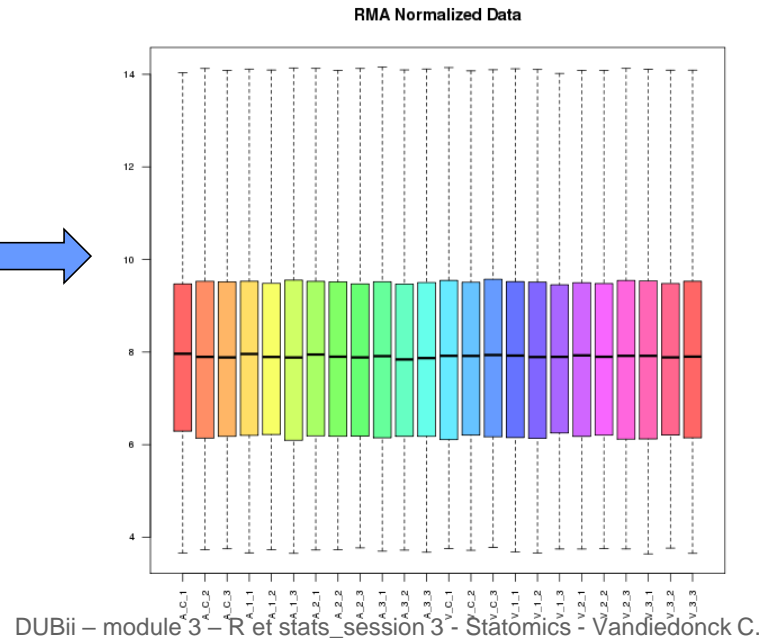
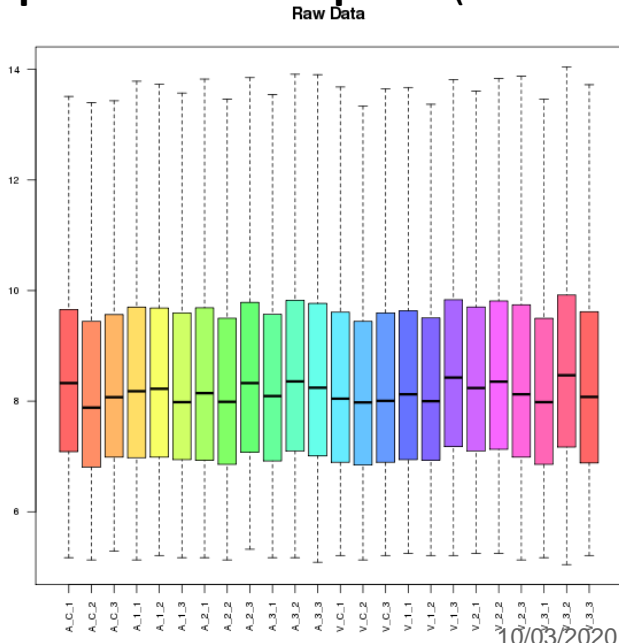
=> *variable quantitative continue*

distribution asymétrique à droite

-> le passage en \log_2 donne souvent une distribution 'normale' (1er sens de normalisation)



Il est aussi nécessaire de **normaliser les échantillons entre eux (2ème sens de normalisation) pour pouvoir les comparer (même échelle)**

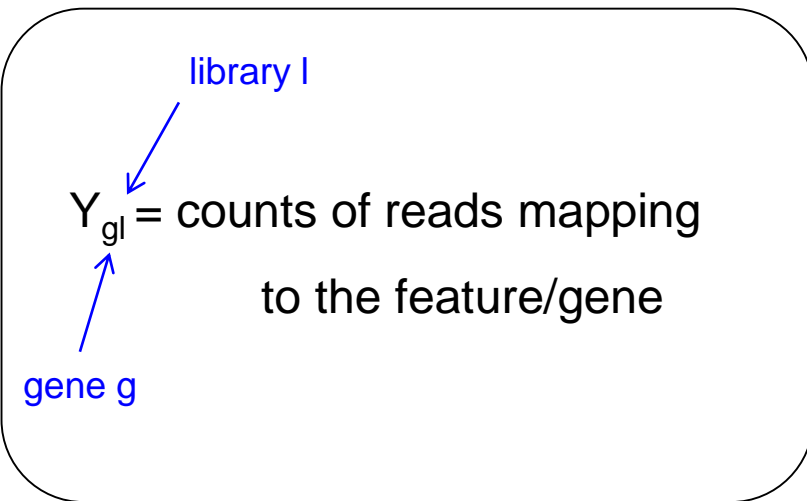


Nature des données d'expression du transcriptome

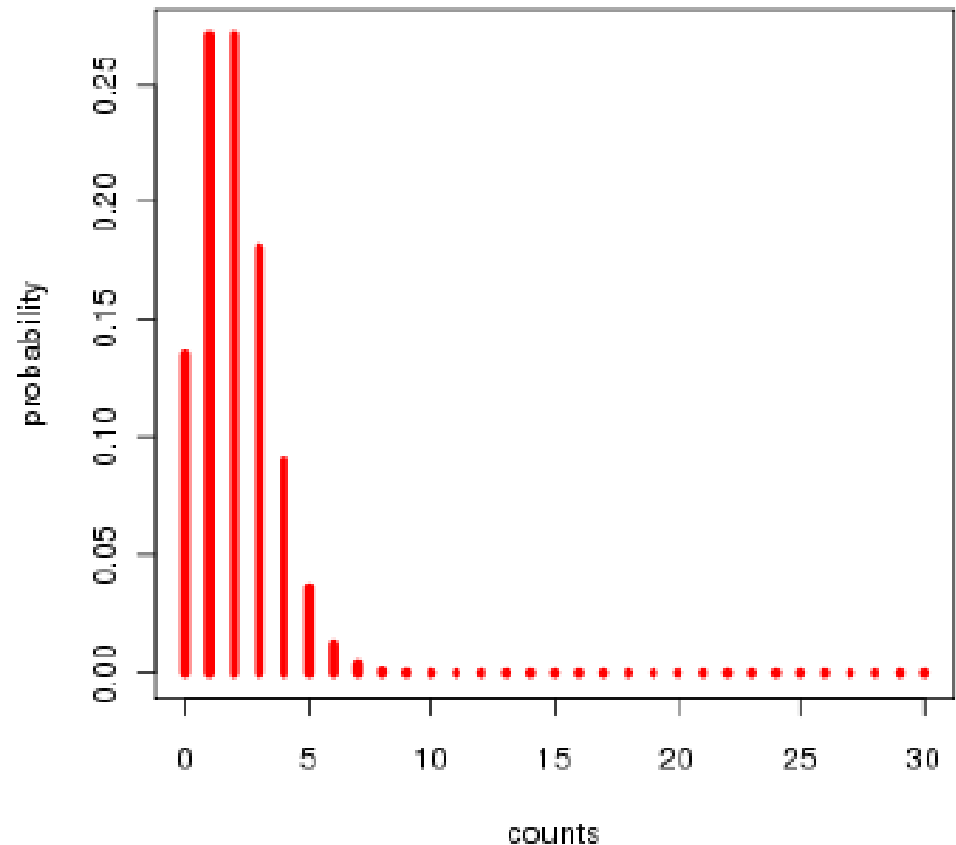
RNASeq

L'abondance des transcrits est mesurée par le nombre de lectures cartographiées au niveau de la séquence génomique du transcrit
= *comptes de lectures*

↪ *Variable quantitative discrète*



Distribution de comptes de RNASeq



=> Il faut utiliser la bonne loi de distribution (Poisson, Négative Binoimale...)

Estimating mean and variance in microarray experiments

Gene expression values are given by fluorescence intensities

- continuous variables
- assumed to fit a Student t distribution (after log2 transformation) of the difference mean

$$t_{\text{gene } i} = \frac{\bar{x}_i}{\tilde{s}_i / \sqrt{n}}$$

- but low number of replicates => difficult to estimate the variance

⇒ LIMMA (Linear Model for MicroArray experiments)

- uses a “moderated” t statistics using information from all genes (group of genes g like gene i) to estimate the variance

$$\tilde{t}_{\text{gene } i} = \frac{\bar{M}_i}{\tilde{s}_g / \sqrt{n}}$$

- allows for linear models
- design matrix => the factors to be accounted for in the model
- contrast matrix => which comparisons are of interest
- accounts for multiple testing: computes adjusted p-value (FDR B-H)

Estimating mean and variance in RNASeq experiments

In RNA-Seq, each feature (gene, exon, isoform) has an expression rate: each segment is sequenced with a low probability

Number of reads from gene g in library i can be captured by a Poisson model (Marioni et al. 2008)

$$r_{ij} \sim \text{Poisson} (\lambda_{ig} = \mu_{ig} k_{ig})$$

where

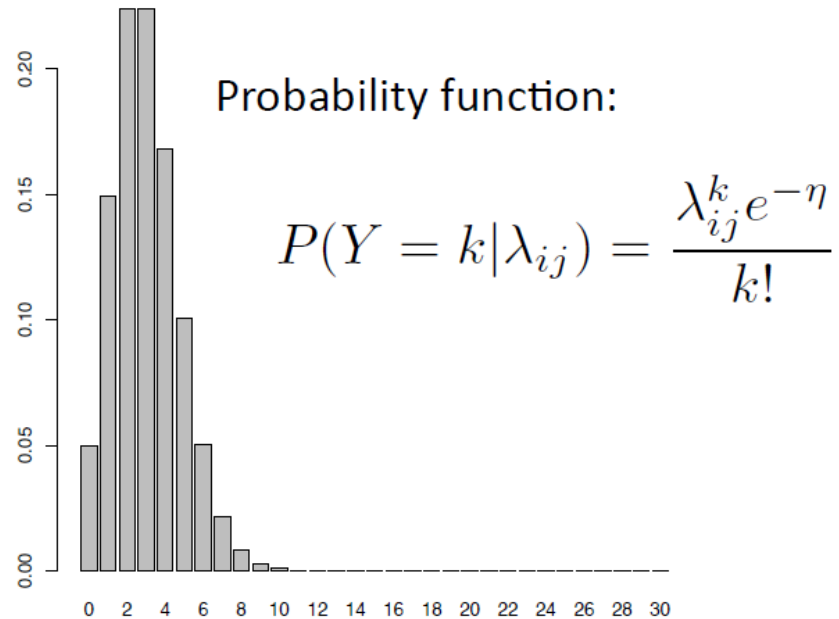
μ_{ig} is the concentration of the RNA

k_{ig} is a normalisation constant

$$\hat{\mu}_{ig} = \frac{r_{ig}}{k_{ig}}$$

$$\lambda_{ig} = \mu_{ig} k_{ig} = E(r_{ij}) = \text{Var}(r_{ij})$$

↳ If n Xiid $\sim \text{Poisson}(\lambda)$, $\sum X_i \sim \text{Poisson}(n\lambda)$

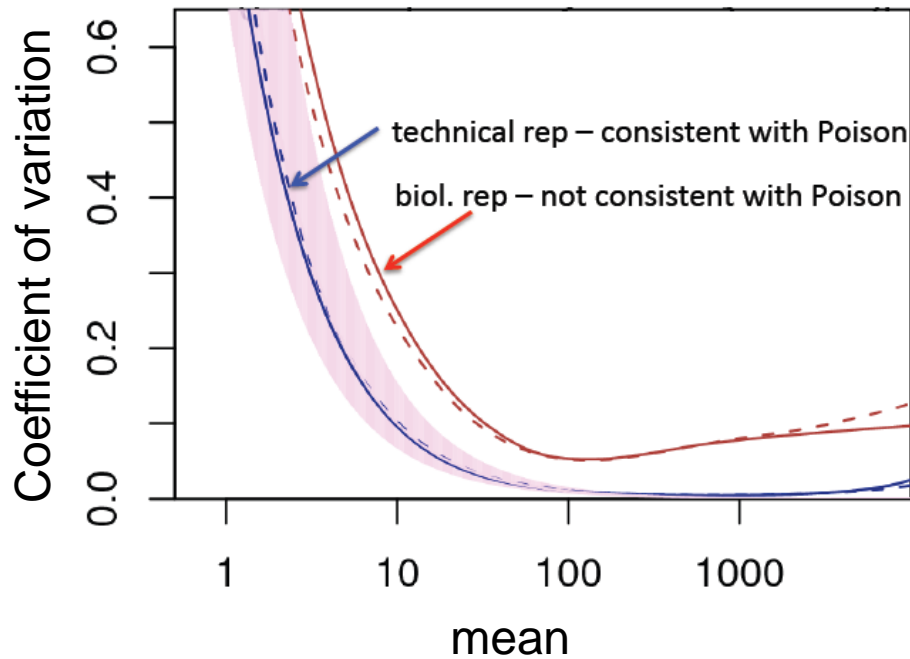


Need to account for extra variability

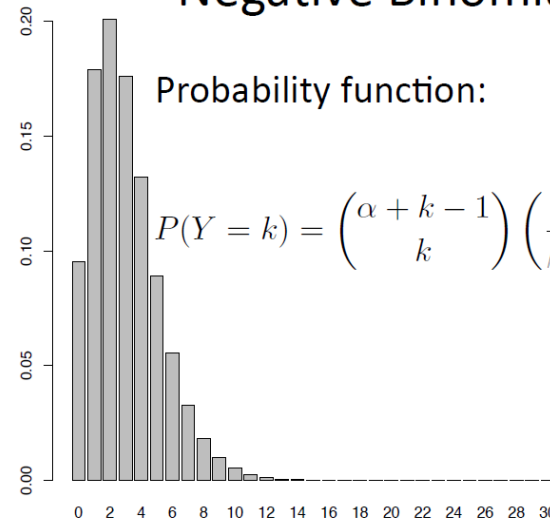
Poisson distribution accounts for technical variation

But biological noise induces an overdispersion

Convergence on a **negative binomial model** for count data



Negative Binomial



Probability function:

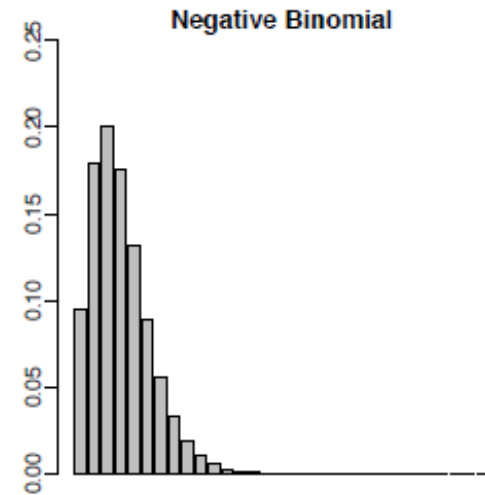
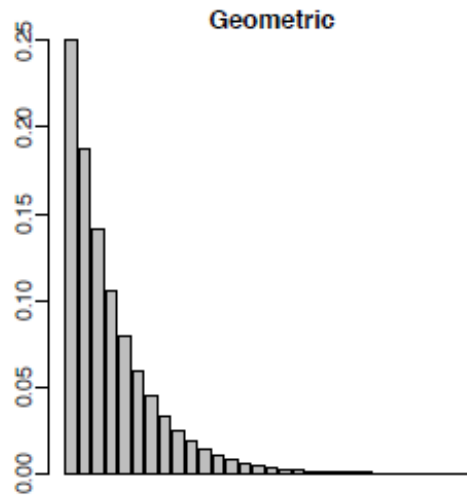
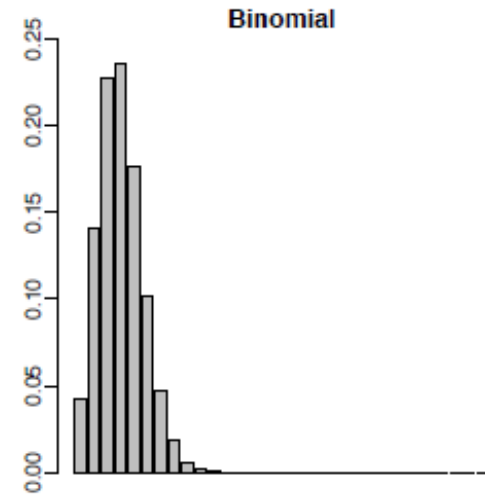
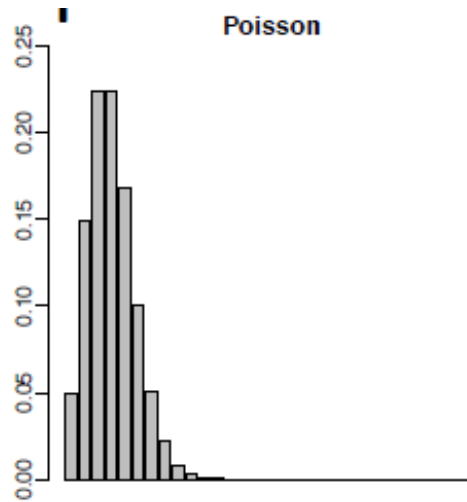
$$P(Y = k) = \binom{\alpha + k - 1}{k} \left(\frac{1}{\beta + 1}\right)^\alpha \left(\frac{\beta}{\beta + 1}\right)^k$$

$$r_{ij} \sim \text{NB}\left(\alpha, \frac{1}{1 + \beta}\right)$$

where α and β are the parameters of a gamma distribution followed by the rates of different samples

Examples of discrete distributions

- **Binomial:** probability of k successes of a Bernoulli variable
- **Geometric:** probability of k failures before 1st success
- **Poisson:** probability of k rare events
- **Negative-binomial:** probability of k failures before n successes



Modelling the variation

The example of DESeq and EdgeR

- generalized linear model fitting the negative binomial distribution:

$$K_{ij} \sim \text{NB}(\mu_{ij}, \alpha_i)$$

K_{ij} : counts of reads for gene i in sample j

α_i : gene-specific dispersion parameter

μ_{ij} : fitted mean

➤ $\mu_{ij} = s_j q_{ij}$

s_j : sample-specific size parameter

q_{ij} : a parameter proportional to the expected true concentration of fragments for sample j

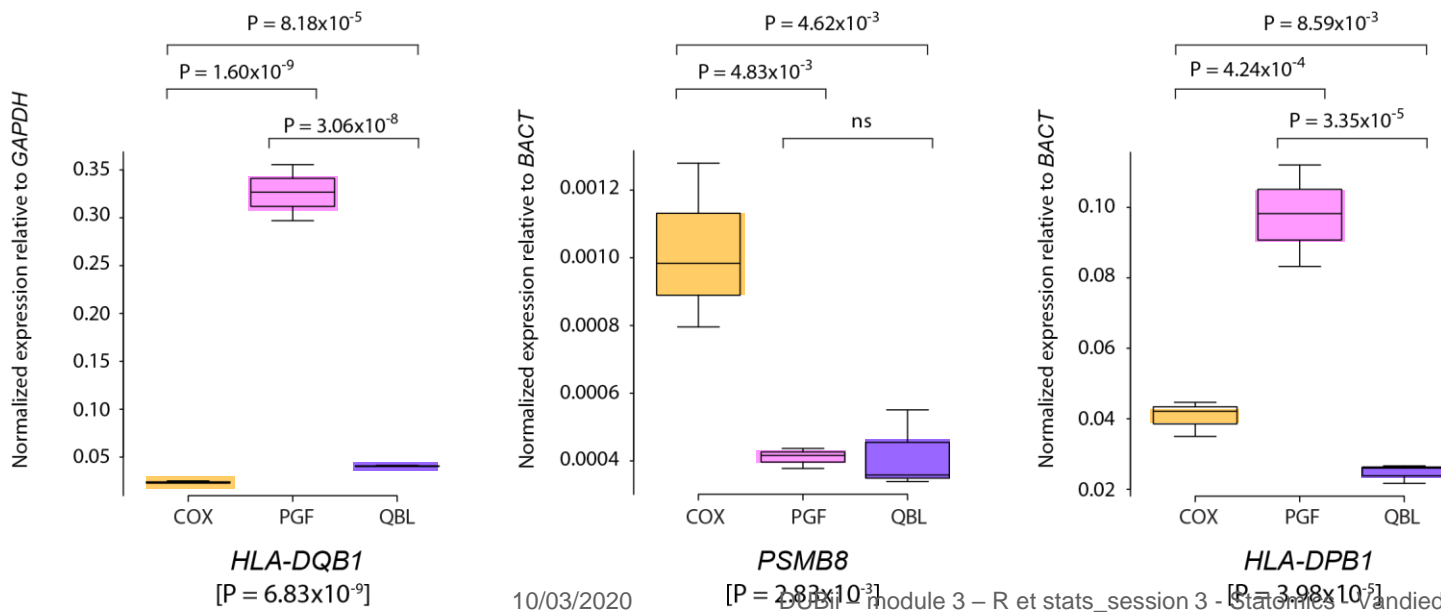
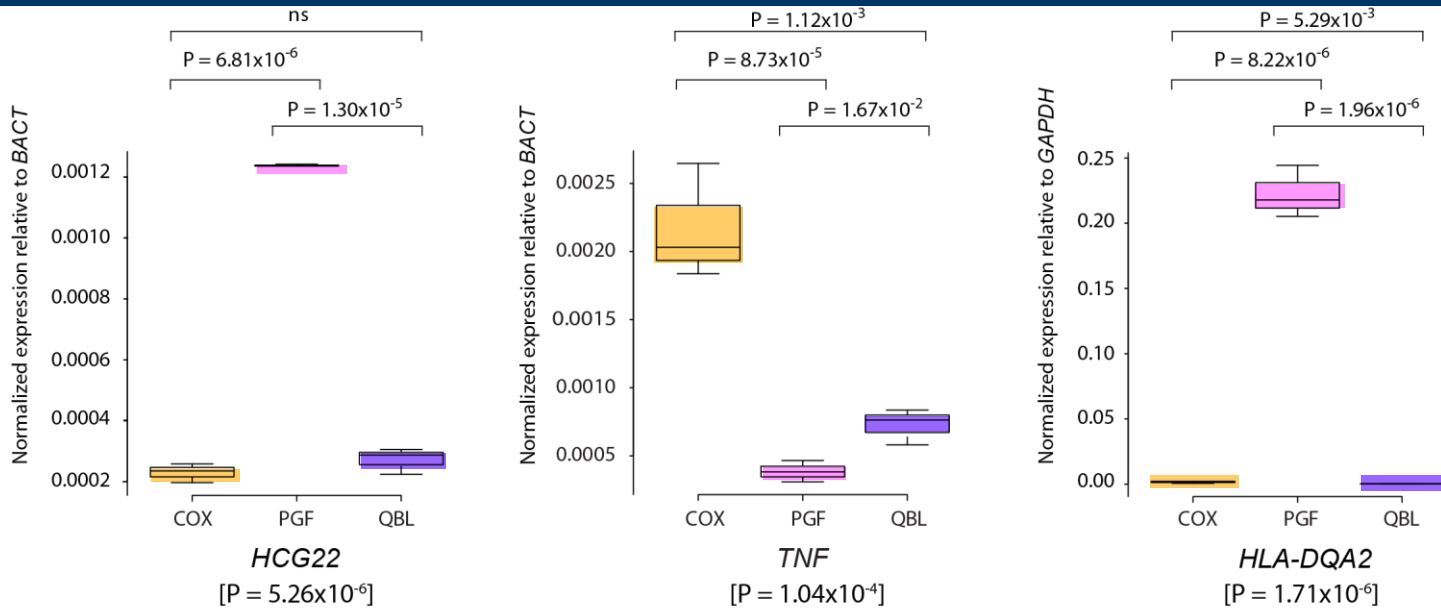
➤ $\log_2(q_{ij}) = x_j \beta_i$

β_i : the log2 fold change for gene i for each column (j .) of the model matrix X

Exemple d'analyse de l'expression différentielle

| Gene Name | Class | log2 (Fold Change) | | | Adj.P.Val |
|-------------------|-------|--------------------|------------|------------|------------------------|
| | | COX vs PGF | QBL vs PGF | QBL vs COX | |
| ZFP57 | I | 2.77 | 0.00 | -2.76 | 1.22x10 ⁻¹⁴ |
| HLA-DPB2 * | II | -3.19 | -3.02 | 0.17 | 2.89x10 ⁻¹² |
| HLA-DQA2 | II | -2.45 | -1.62 | 0.82 | 1.91x10 ⁻¹¹ |
| HLA-DQB2 | II | -2.74 | -2.58 | 0.16 | 3.21x10 ⁻¹¹ |
| HLA-21 * | I | -2.52 | 0.36 | 2.87 | 1.32x10 ⁻¹⁰ |
| TNF | III | 1.90 | 1.03 | -0.87 | 4.79x10 ⁻¹⁰ |
| HLA-DPB1 | II | -2.08 | -0.90 | 1.18 | 6.44x10 ⁻¹⁰ |
| RPL32P1 * | II | -1.52 | -1.19 | 0.33 | 2.07x10 ⁻⁰⁹ |
| HLA-B | I | -0.06 | -1.19 | -1.13 | 6.59x10 ⁻⁰⁹ |
| HLA-A | I | -1.51 | -1.86 | -0.35 | 2.30x10 ⁻⁰⁸ |
| HLA-L * | I | -1.29 | -1.47 | -0.18 | 2.30x10 ⁻⁰⁸ |
| XXbac-BPG254F23.6 | II | -1.59 | -1.59 | 0.00 | 2.50x10 ⁻⁰⁸ |
| HCG22 | I | -1.56 | -1.26 | 0.30 | 2.96x10 ⁻⁰⁸ |
| XXbac-BPG254F23.5 | II | -1.42 | -1.61 | -0.19 | 1.33x10 ⁻⁰⁷ |
| LTA | III | 1.32 | 0.57 | -0.75 | 2.04x10 ⁻⁰⁷ |
| NCR3 | III | 0.87 | 0.95 | 0.08 | 4.95x10 ⁻⁰⁷ |
| HLA-F | I | 0.15 | -0.90 | -1.05 | 4.95x10 ⁻⁰⁷ |
| HLA-DOA | II | -1.32 | -0.89 | 0.43 | 5.07x10 ⁻⁰⁷ |
| TAP1 | II | 0.97 | 0.08 | -0.89 | 6.86x10 ⁻⁰⁷ |
| LTB | III | -0.95 | -0.06 | 0.89 | 7.02x10 ⁻⁰⁷ |
| LST1 | III | -0.18 | 0.48 | 0.66 | 9.42x10 ⁻⁰⁷ |
| DAQB-335A13.8 | I | 0.61 | -0.02 | -0.63 | 1.12x10 ⁻⁰⁶ |
| TCF19 | I | 1.11 | 0.62 | -0.49 | 1.49x10 ⁻⁰⁶ |
| CLIC1 | III | 1.22 | 0.57 | -0.66 | 1.49x10 ⁻⁰⁶ |
| HLA-DMA | II | -0.57 | -0.89 | -0.33 | 3.52x10 ⁻⁰⁶ |
| BRD2 | II | 0.78 | 0.27 | -0.51 | 3.60x10 ⁻⁰⁶ |
| NRM | I | 0.77 | 0.39 | -0.38 | 4.48x10 ⁻⁰⁶ |
| HLA-C | I | 0.05 | 1.11 | 1.06 | 4.98x10 ⁻⁰⁶ |
| PSMB9 | II | 0.42 | -0.29 | -0.71 | 6.05x10 ⁻⁰⁶ |
| HCG27 | I | 0.56 | 0.06 | -0.50 | 7.01x10 ⁻⁰⁶ |

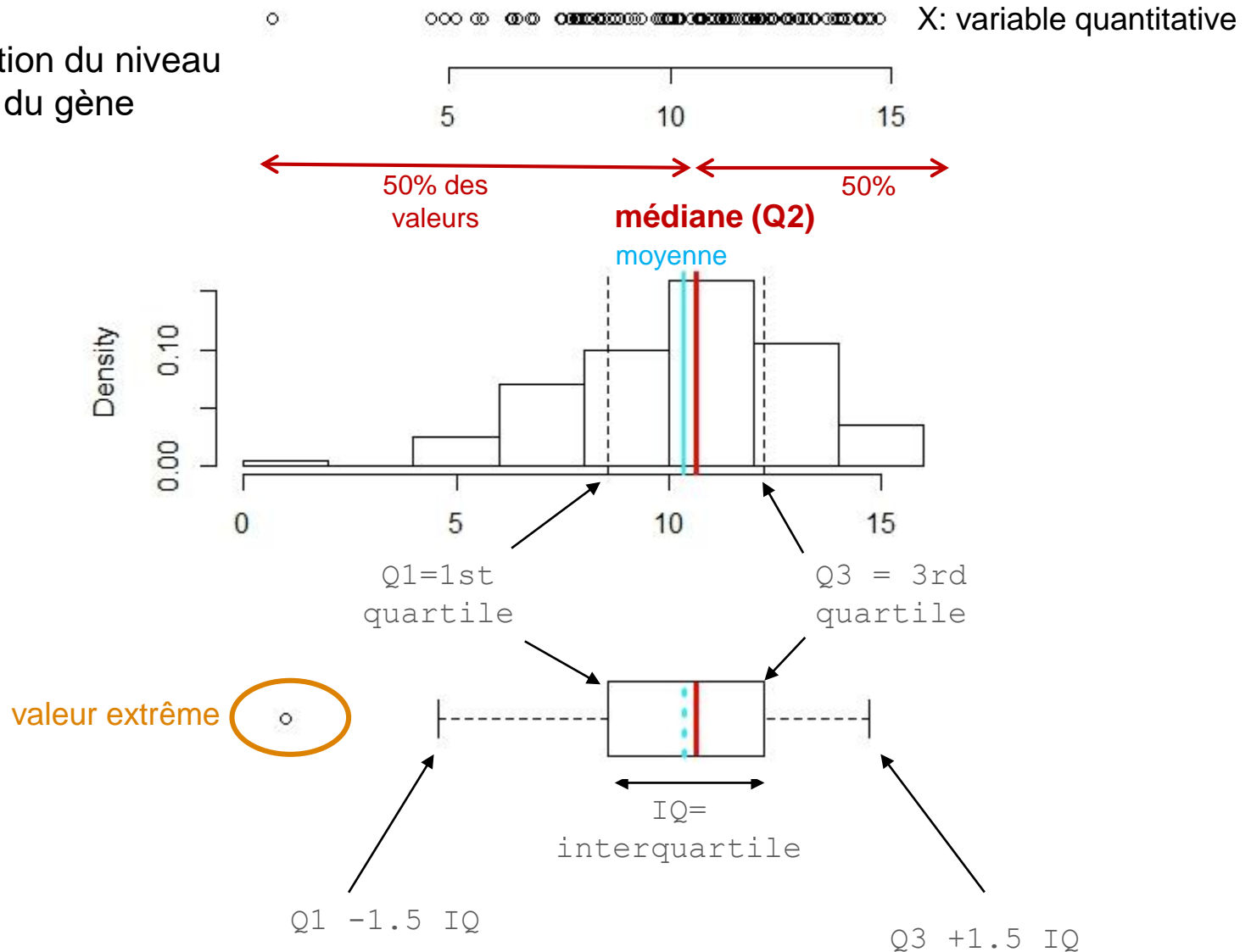
Validation des meilleurs gènes par qPCR



Représentation graphique de l'analyse différentielle

Boxplot:

de la distribution du niveau d'expression du gène



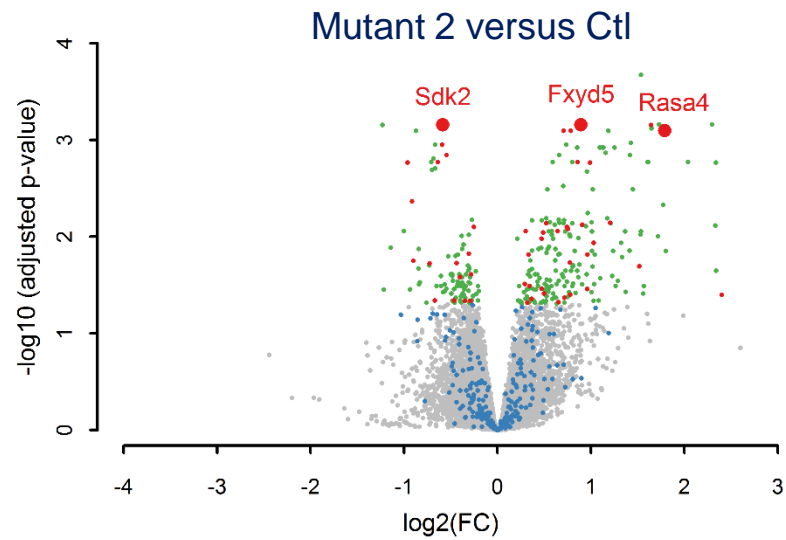
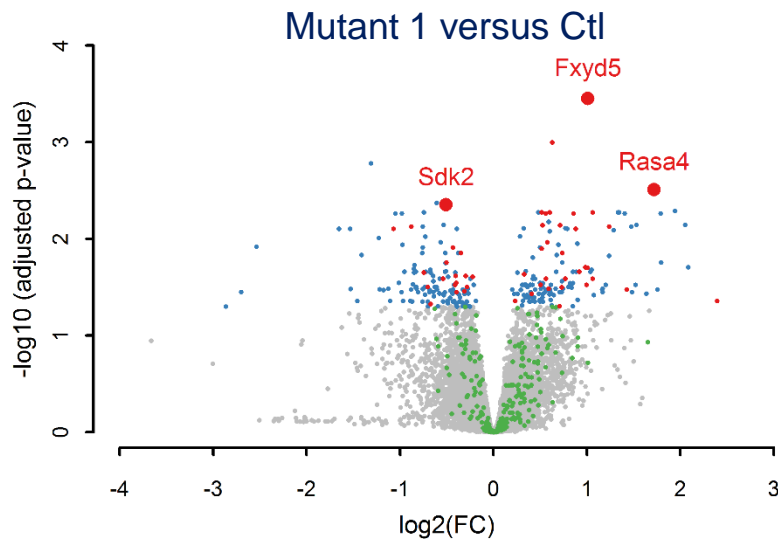
Représentation graphique de l'analyse différentielle

Volcano plots:

$X = \log_2(\text{Fold change})$

$Y = -\log_{10}(\text{pvalue})$

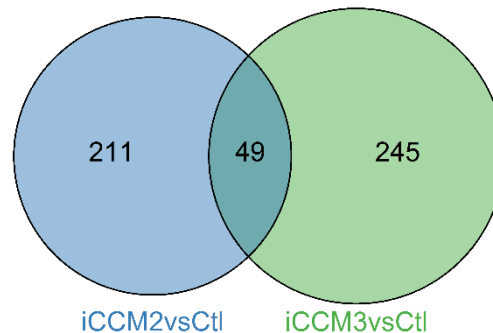
➤ Exemple ici chez la souris avec 2 gènes KO versus Wild Type



● not significant ● mutant 1-specific ● mutant 2-specific ● shared

Diagramme de Venn

intersection des listes de gènes



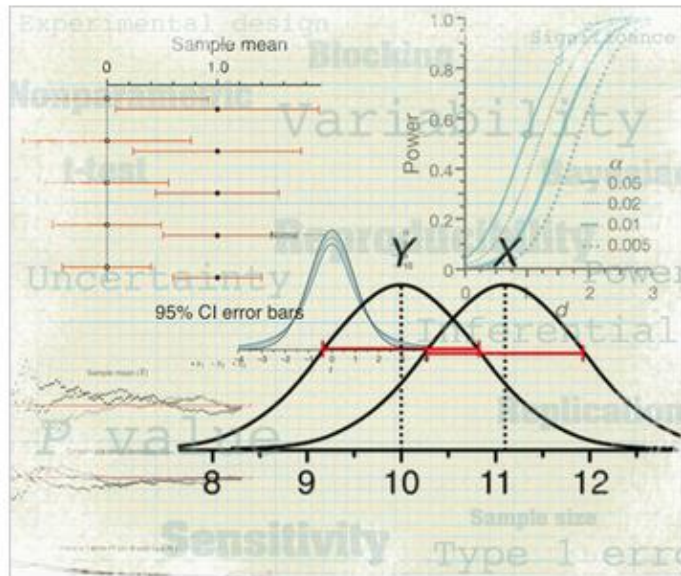
4. The 3rd issue: reducing dimensionality

-> cf. next sessions

5. Liens

Statistics for biologists

Home | [Practical guides](#) | [Statistics in biology](#) | [Points of Significance](#) | [Other resources](#)



There is no disputing the importance of statistical analysis in biological research, but too often it is considered only after an experiment is completed, when it may be too late.

This collection highlights important statistical issues that biologists should be aware of and provides practical advice to help them improve the rigor of their work.

Nature Methods' **Points of Significance** column on statistics explains many key statistical and experimental design concepts. **Other resources** include an online plotting tool and links to statistics guides from other publishers.

Image Credit: Erin DeWalt

Statistics in biology

Nature News | Editorial

Number crunch



Top picks
from **nature** news

Nature News | News

Points of significance: <http://mkweb.bcgsc.ca/pointsofsignificance/>

RESOURCES

@ Martin Krzywinski / Genome Sciences Center / mkweb.bcgsc.ca on Twitter Lips that taste of tears, they say, are the best for kissing. • Dorothy Parker • get cranky

ART IS SCIENCE IS ART

Share

mkweb.bcgsc.ca/pointsofsignificance/bibliography.mhtml

STATISTICS EXPLAINED

POINTS OF VIEW

VIZBI 2012 VISUALIZATION PRINCIPLES TUTORIAL

ADOBE SWATCHES FOR BREWER PALETTES

PRESENTATIONS

Let me tell you about something.

DRINKS & SCIENCE WORKSHOP

VIZBI 2013 KEYNOTE VISUAL DESIGN PRINCIPLES

SCIENCE CENTRE

FRUITIGER ABCDEFGHIJKLMNOPQRSTUVWXYZAAO

BC Cancer Agency CARE & RESEARCH An agency of the Provincial Health Services Authority

Circos is back for 4rd year at 2014 Bioinformatics and Comparative Genome Analysis course by the Pasteur Institute—Athens May 7

THINGS ON THE SIDE

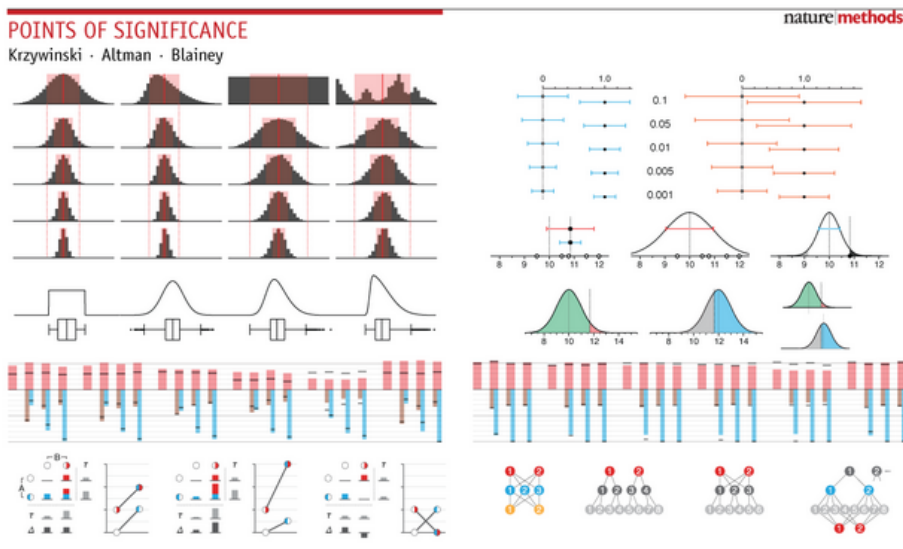
quotes · typography · art of π · road trips · hitchhens · ascii · lit 2.0 · questions · writing · words · satire · lost · photo · universe · keyboards · time · lexical · covers · choices · clocks · rockets · famous rat · lotro

STATISTICS + DATA

NATURE METHODS: POINTS OF SIGNIFICANCE

POINTS OF SIGNIFICANCE

Krzywinski · Altman · Blainey



nature methods

NEWS + THOUGHTS

POINTS OF SIGNIFICANCE COLUMN NOW OPEN ACCESS

Tue 10-02-2015

Nature Methods has announced the launch of a new statistics collection for biologists.

Altman, N. & Krzywinski, M. (2015) [Points of Significance: Sources of Variation](#) *Nature Methods* 12:5-6.

Krzywinski, M., Altman, N. & Blainey, P. (2014) [Points of Significance: Two factor designs](#) *Nature Methods* 11:1187-1188.

Krzywinski, M., Altman, N. & Blainey, P. (2014) [Points of Significance: Nested designs](#) *Nature Methods* 11:977-978.

Blainey, P., Krzywinski, M. & Altman, N. (2014) [Points of Significance: Replication](#) *Nature Methods* 11:879-880.

Krzywinski, M. & Altman, N. (2014) [Points of Significance: Analysis of variance \(ANOVA\) and blocking](#) *Nature Methods* 11:699-700.

Krzywinski, M. & Altman, N. (2014) [Points of Significance: Designing comparative experiments](#) *Nature Methods* 11:597-598.

Krzywinski, M. & Altman, N. (2014) [Points of Significance: Non parametric tests](#) *Nature Methods* 11:467-468.

Krzywinski, M. & Altman, N. (2014) [Points of Significance: Comparing samples—Part II — Multiple Testing](#) *Nature Methods* 11:355-356.

Krzywinski, M. & Altman, N. (2014) [Points of Significance: Comparing samples—Part I — t-tests](#) *Nature Methods* 11:215-216.

Krzywinski, M. & Altman, N. (2014) [Points of Significance: Visualizing samples with box plots](#) *Nature Methods* 11:119-120.

Krzywinski, M. & Altman, N. (2013) [Points of Significance: Power and sample size](#) *Nature Methods* 10:1139-1140.

Krzywinski, M. & Altman, N. (2013) [Points of Significance: Significance, P values and t-tests](#) *Nature Methods* 10:1041-1042.

Krzywinski, M. & Altman, N. (2013) [Points of Significance: Error bars](#) *Nature Methods* 10:921-922.

Krzywinski, M. & Altman, N. (2013) [Points of Significance: Importance of being uncertain](#) *Nature Methods* 10:809-810.

▲ Points of Significance column in Nature Methods. (Launch of Points of Significance)

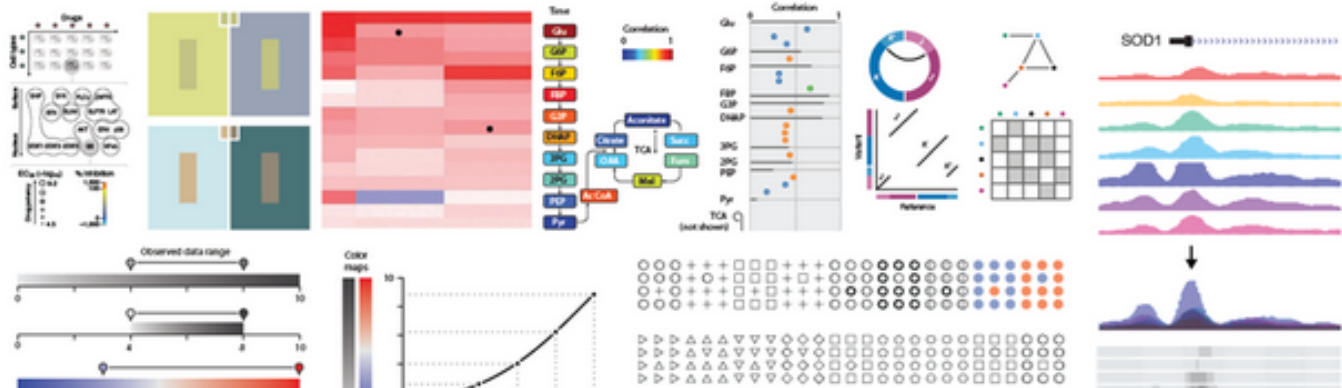
Points of view: <http://mkweb.bcgsc.ca/pointsofview>

COMMUNICATION + SCIENCE

NATURE METHODS: POINTS OF VIEW

POINTS OF VIEW

Wong · Krzywinski · Gehlenborg · Nielsen · Soresh · Kjægaard · Savig · Cairo



FULL COLLECTION

open access through August 2013

▲ The full collection of a 35 Points of View column is now available. (3 years of Points of View)

PRACTICAL TIPS FOR EFFECTIVE FIGURES

POINTS OF VIEW – HISTORY

In its 2.5 year history, the PoV column has established a significant legacy— it is one of the most frequently accessed parts of Nature Methods. The reason I think is clear: the community sees the value in clear and effective visual communication and acknowledges the need for a forum in which best practices in the field are presented practically and accessibly.

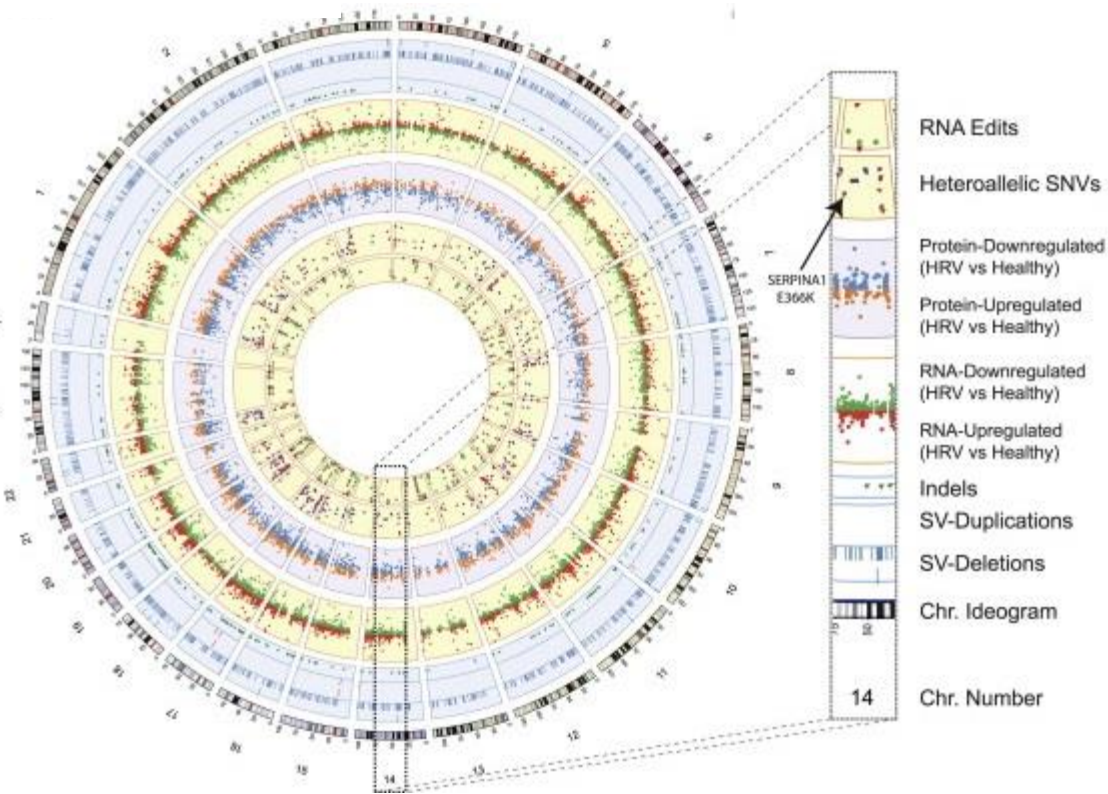
Bang Wong, in collaboration with visiting authors (Noam Shores, Nils Gehlenborg, Cydney Nielsen and Rikke Schmidt Kjægaard), has penned 29 columns in the period of August 2010 to December 2012, covering broad topics such as salience, Gestalt principles, color, typography, negative space, layout, and data integration.

When it was A.C. Greyling's turn to speak at a debate in which Christopher Hitchens and Richard Dawkins already made their points, Greyling said

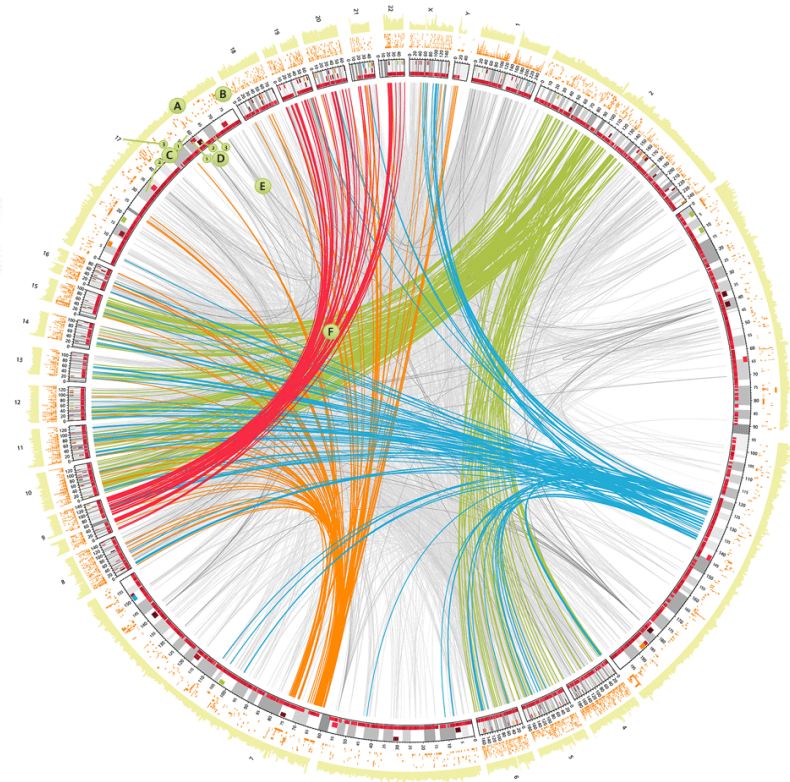
Circos to represent genomic traits:

http://circos.ca/intro/genomic_data/

Co-localisation



Interaction



Cell 148, 1293–1307, March 16, 2012

Personal Omics Profiling Reveals Dynamic Molecular and Medical Phenotypes

Rui Chen,^{1,11} George I. Mias,^{1,11} Jennifer Li-Pook-Than,^{1,11} Lihua Jiang,^{1,11} Hugo Y.K. Lam,^{1,12} Rong Chen,^{2,12} Elana Miriami,¹ Konrad J. Karczewski,¹ Manoj Hariharan,¹ Frederick E. Dewey,³ Yong Cheng,¹ Michael J. Clark,¹ Hogune Im,¹ Lukas Habegger,^{6,7} Suganthi Balasubramanian,^{6,7} Maeve O'Huallachain,¹ Joel T. Dudley,² Sara Hillenmeyer,¹ Rajini Haraksingh,¹ Donald Sharon,¹ Ghia Euskirchen,¹ Phil Lacroute,¹ Keith Bettinger,¹ Alan P. Boyle,¹ Maya Kasowski,¹ Fabian Grubert,¹ Scott Seki,² Marco Garcia,² Michelle Whirl-Carrillo,¹ Mercedes Gallardo,^{9,10} Maria A. Blasco,⁹ Peter L. Greenberg,⁴ Phyllis Snyder,¹ Teri E. Klein,¹ Russ B. Altman,^{1,5} Atul J. Butte,² Euan A. Ashley,³ Mark Gerstein,^{6,7,8} Kari C. Nadeau,² Hua Tang,¹ and Michael Snyder^{1,*}

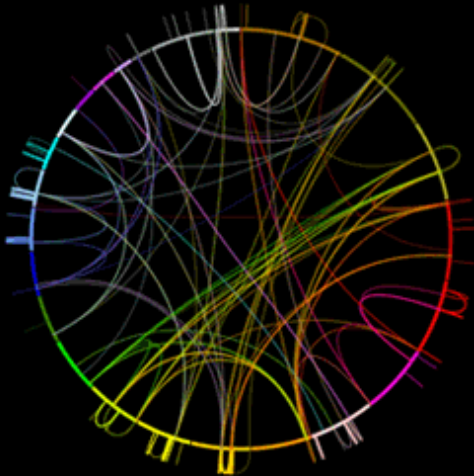
10/03/2020

DUBii – module 3 – R et stats_session 3 - Statomics - Vandiedonck C.

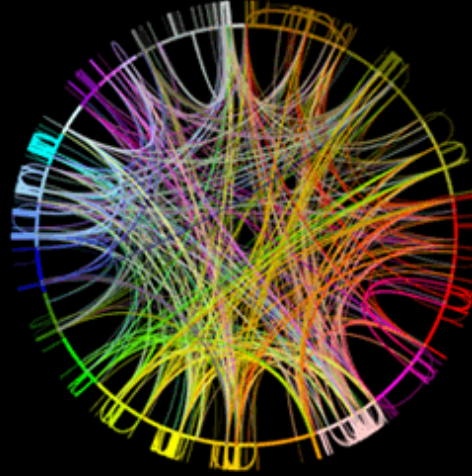
47 / 48

Towards an increasing complexity of omics!

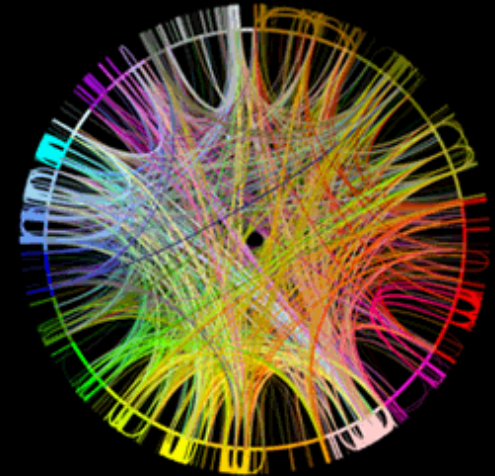
GENOME



OMEOME



METAOMEOME



And rather quickly
it has come to this.