



Sciences Economiques & Sociales de la Santé  
& Traitement de l'Information Médicale

[www.sesstim-orspaca.org](http://www.sesstim-orspaca.org)

**Laura RICHERT**

*MD, PhD*

*INSERM U1219 - INRIA SISTIM - ISPED, Université de Bordeaux*

## **Big data en épidémiologie**

mai 2016



**Cliquez ici pour voir l'intégralité des ressources associées à ce document**

Webinar QuanTIM  
20 mai 2016

# Big data en épidémiologie

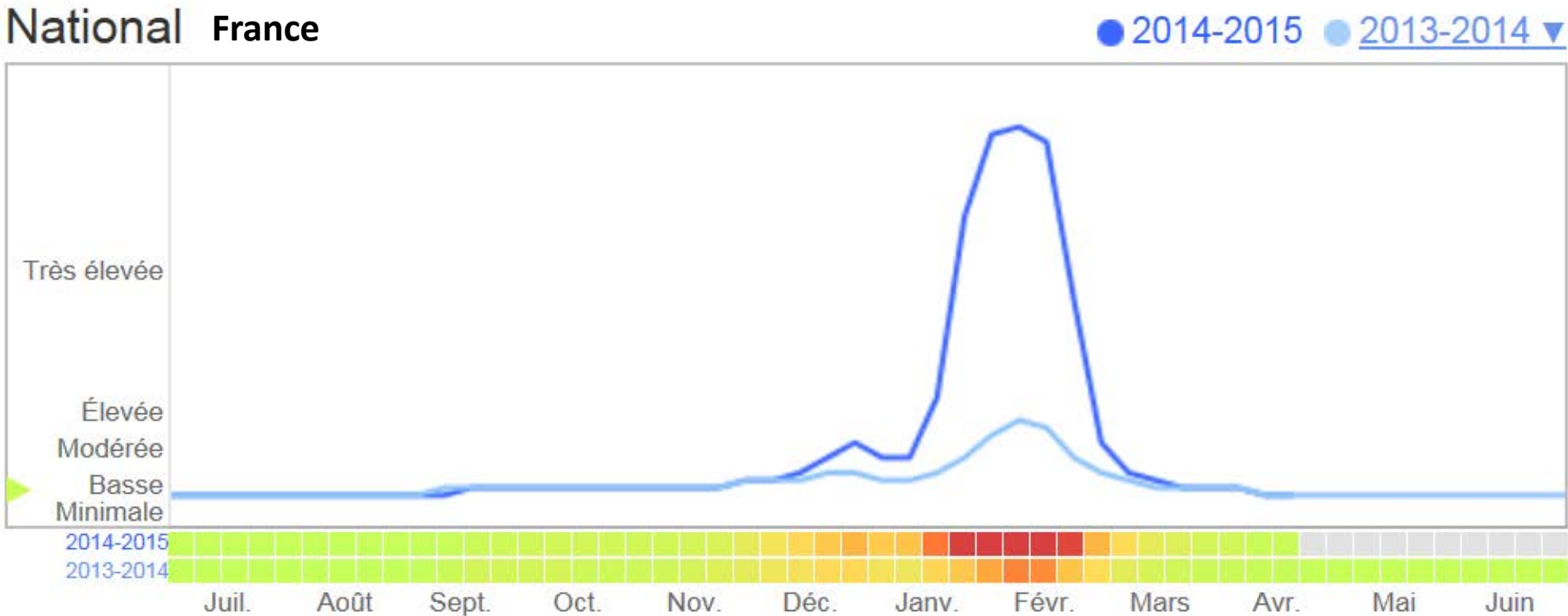
Laura Richert

INSERM U1219/INRIA SISTM  
ISPED, Université de Bordeaux



# Big Data

<https://www.google.org/flutrends/fr/#FR>



OPEN ACCESS Freely available online

PLoS one

## Assessing Google Flu Trends Performance in the United States during the 2009 Influenza Virus A (H1N1) Pandemic

Samantha Cook<sup>1</sup>, Corrie Conrad<sup>2\*</sup>, Ashley L. Fowlkes<sup>3</sup>, Matthew H. Mohebbi<sup>1</sup>

# Séquençage génétique

2005

2010

2012



13 Mb/hour  
\$10/Mb

Pacific Biosciences



200 Mb/hour  
\$2/Mb

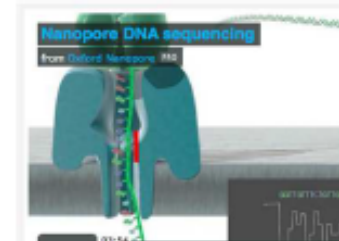


0,5 Gb/hour  
\$1/Mb

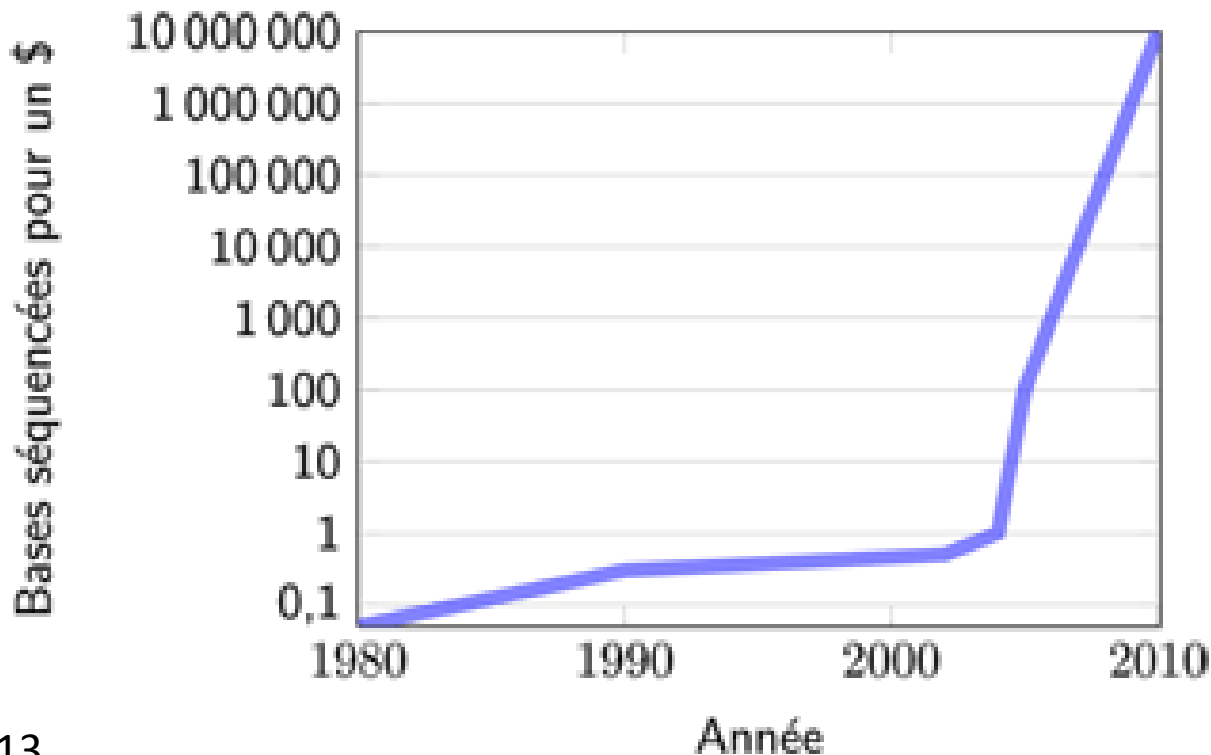
illumina



2,5 Gb/hour  
\$0.1/Mb



nanopore



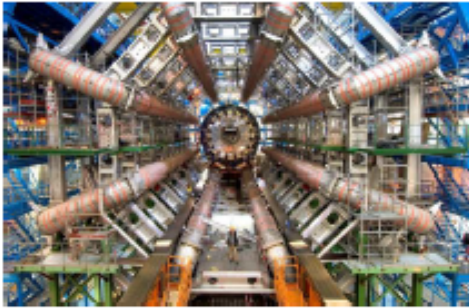
**WHAT DO WE MEAN BY « BIG DATA »?**

**QUE VEUT-ON DIRE PAR « DONNEES  
MASSIVES »?**

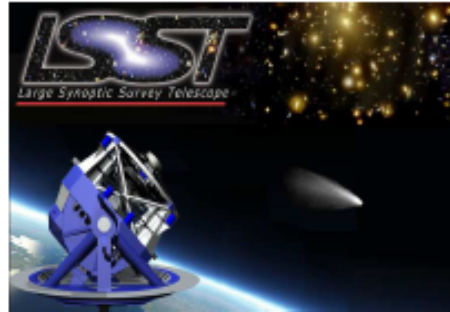
- **Edition 2016 du *Robert***

- Big data **n. m.** (mots anglais « données massives») anglic. Le big data : l'ensemble des données générées par les nouvelles technologies, caractérisées par leur volume colossal. – recomm. offic. *mégadonnées* **n. f. pl.**

Un film de 1h30 (qualité dvd): 2.5 Go



LHC  
~ 1 To/h (raw: 300 Go/h)  
15 Po/year



LSST  
~ 1 To/h (raw: 40 Tb/h)  
100 Po/year

facebook

4 Po / year

You Tube™

1 h vidéo / sec  
(1 année = 36 siècles de TV classique)

1 bit : 0/1  
1 octet = 8 bits  
1 Mo =  $10^6$  octets  
1 Go =  $10^9$  octets  
1 To =  $10^{12}$  octets  
1 Po =  $10^{15}$  octets



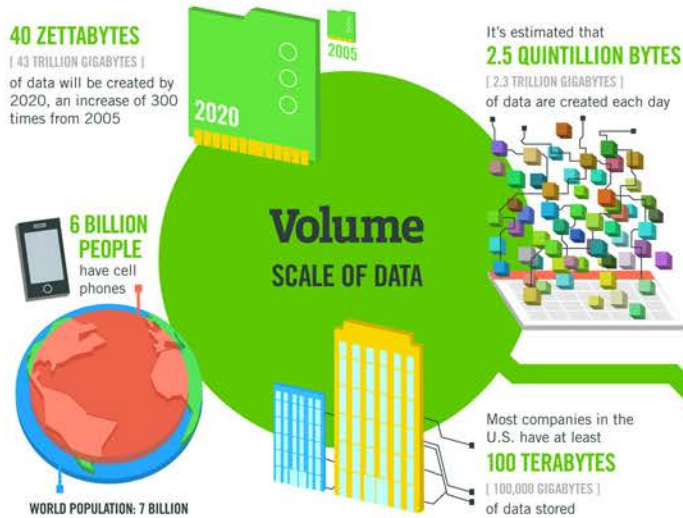
600 Gb/run  
~ 5 Go/h

BGI : 20 Po (sur 8 sites)

FranceGenomique: 5 Po  
IFB: 1 Po



# Les 4 dimensions des big data



## The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015  
**4.4 MILLION IT JOBS**  
will be created globally to support big data, with 1.9 million in the United States



As of 2011, the global size of data in healthcare was estimated to be

**150 EXABYTES**  
[ 161 BILLION GIGABYTES ]



**30 BILLION PIECES OF CONTENT**  
are shared on Facebook every month



By 2014, it's anticipated there will be **420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

**4 BILLION+ HOURS OF VIDEO**  
are watched on YouTube each month



**400 MILLION TWEETS**  
are sent per day by about 200 million monthly active users



## Variety DIFFERENT FORMS OF DATA

The New York Stock Exchange captures **1 TB OF TRADE INFORMATION** during each trading session



Modern cars have close to **100 SENSORS** that monitor items such as fuel level and tire pressure

## Velocity ANALYSIS OF STREAMING DATA

By 2016, it is projected there will be **18.9 BILLION NETWORK CONNECTIONS** - almost 2.5 connections per person on earth



**1 IN 3 BUSINESS LEADERS** don't trust the information they use to make decisions



**27% OF RESPONDENTS**

in one survey were unsure of how much of their data was inaccurate

## Veracity UNCERTAINTY OF DATA

Poor data quality costs the US economy around **\$3.1 TRILLION A YEAR**



# Et en épidémiologie?

From “Big Epidemiology” to “Colossal Epidemiology”  
*When All Eggs Are in One Basket*

*Miguel A. Hernán<sup>a</sup> and David A. Savitz<sup>b</sup>*

Is Size the Next Big Thing in Epidemiology?

*Sengwee Toh and Richard Platt*

*Epidemiology* • Volume 24, Number 3, May 2013

# Epidémiologie classique



Population

Modèle de régression

$$\text{Logit } P(Y=1) = X^T \beta$$



# Genome Wide Association Study



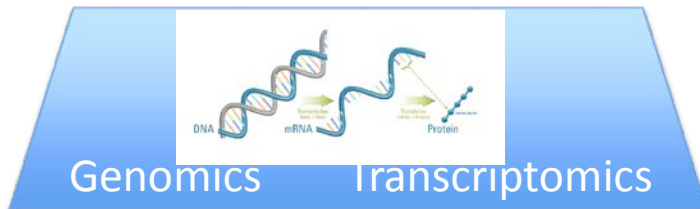
Population

Modèle de régression

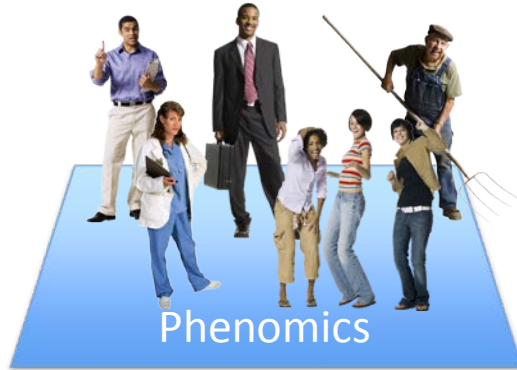
$$\text{Logit } P(Y=1) = X^T \beta$$



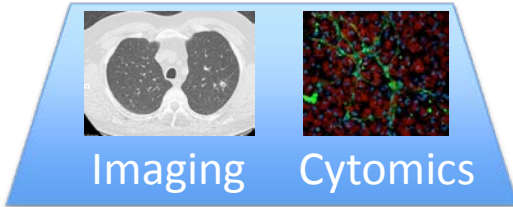
Gènes



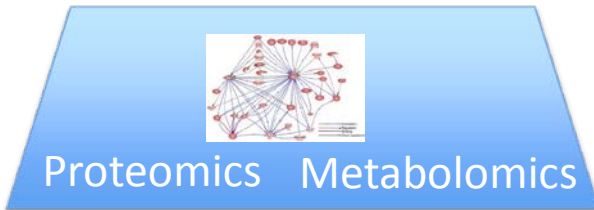
# Epidémiologie systémique



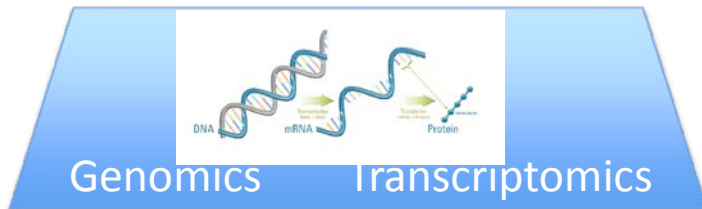
Population



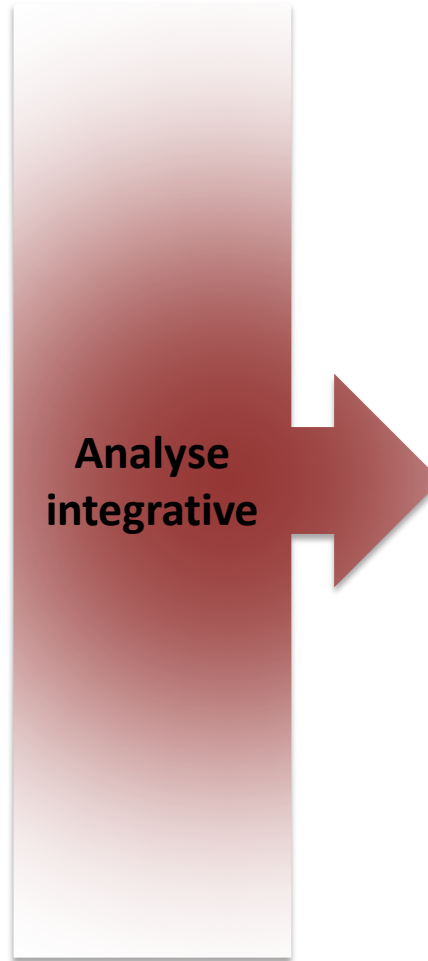
Tissus  
Cellules



Molécules

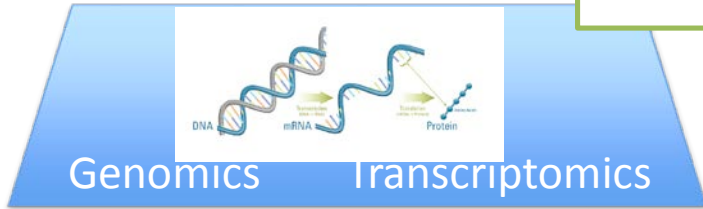
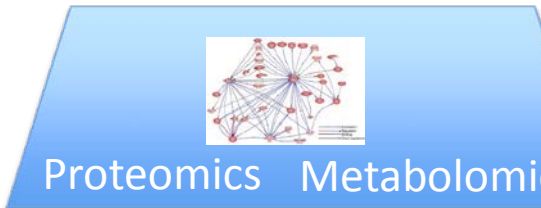
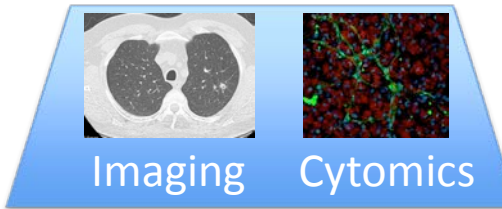


Gènes



Qui?  
Comment?

# Big data



e.g. health national insurance database for 57 M people

## Assurance Maladie



- 1,2 milliards de feuilles de soins gérées chaque année
- Capacité de stockage égale à 450 Téraoctets
- 7 dictionnaires avec 785 000 objets (tables, index, synonymes ...)
- 3 bases de données d'une volumétrie supérieure à 25 téraoctets



# **PB N°1: LA RÉPÉTITION DES TESTS (NOMBREUSES VARIABLES)**

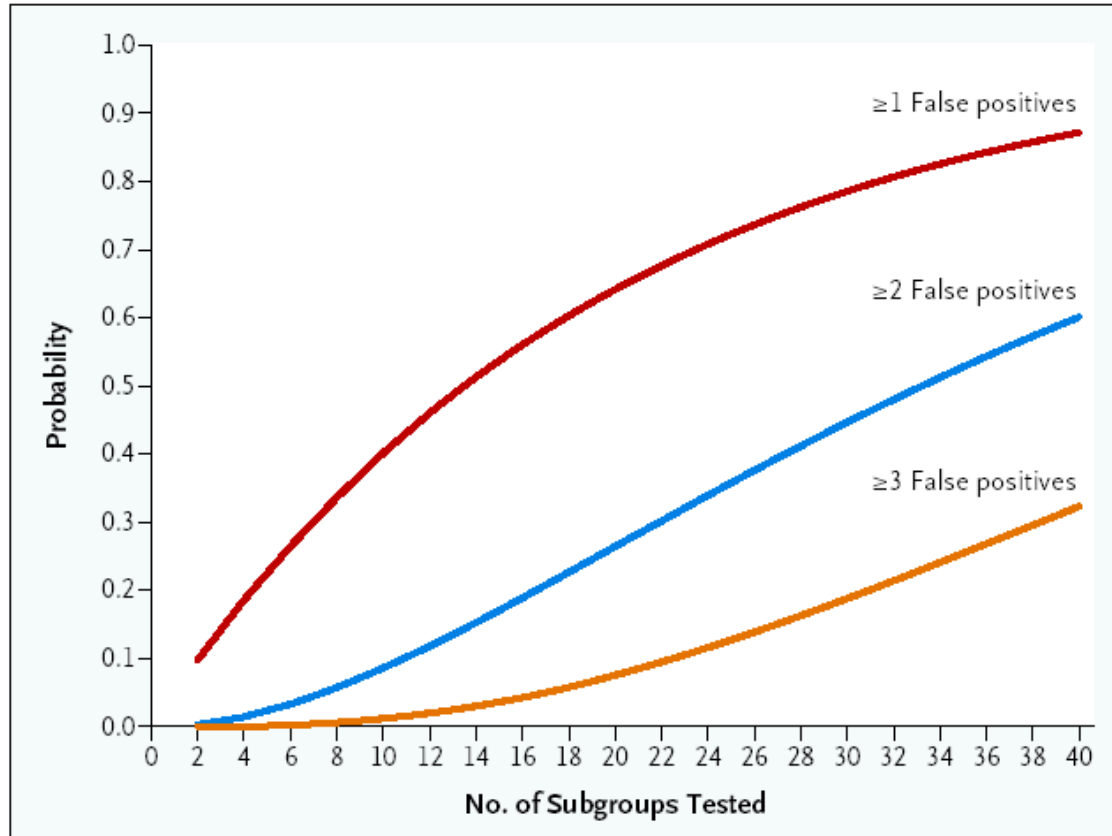
# Multiplicité des tests statistiques

STATISTICS AND MEDICINE

## The Challenge of Subgroup Analyses — Reporting without Distorting

Stephen W. Lagakos, Ph.D.

Related article, page 1706



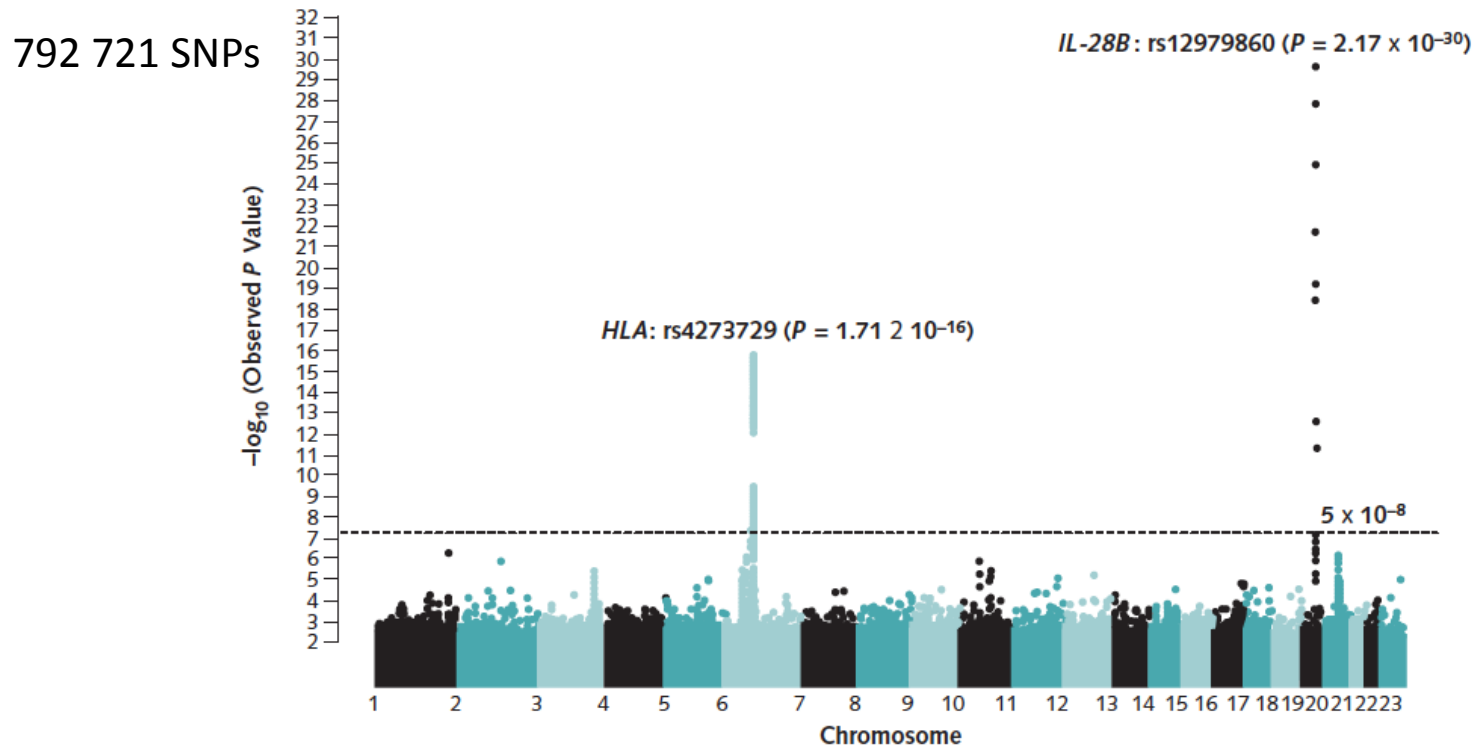
Probability That Multiple Subgroup Analyses Will Yield at Least One (Red), Two (Blue), or Three (Yellow) False Positive Results.



# Etude d'associations génétiques pangénomiques GWAS (genome-wide association study)

« Manhattan plot »

Figure 1. Manhattan plot summarizing the genome-wide association results in 919 persons with spontaneous resolution of hepatitis C virus infection and 1482 persons with chronic hepatitis C virus infection.



Each point corresponds to a  $P$  value from a test of association for a single nucleotide polymorphism. The  $-\log_{10} P$  values are plotted by location of the person's single nucleotide polymorphism across the genome. The dashed line represents an accepted level of genome-wide significance ( $P = 5 \times 10^{-8}$ ). Single nucleotide polymorphisms in the *HLA* and *IL-28B* region on chromosomes 6 and 19, respectively, exceed this threshold. *IL-28B* = *interleukin-28B*.

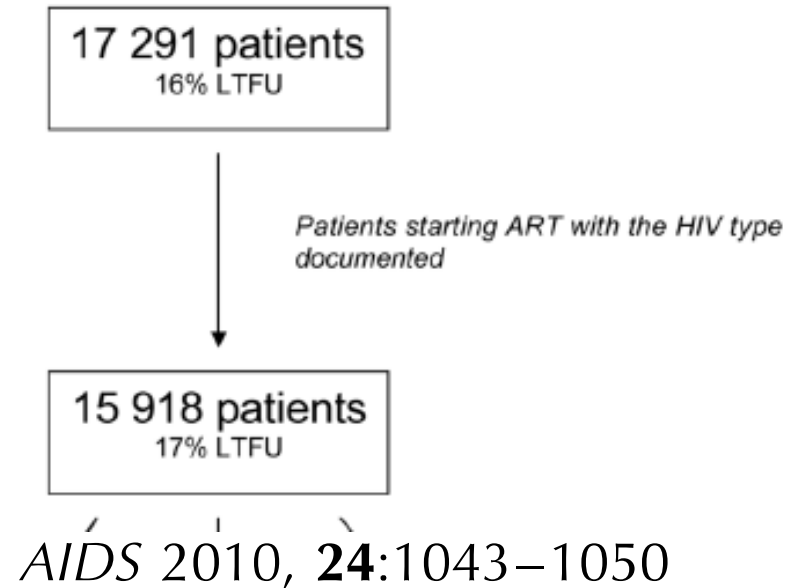
# **PB N°2: LES DONNÉES MANQUANTES**

- Facteurs de confusion connus non mesurés
- Sorties d'étude, perdus de vue



International Epidemiologic Databases to Evaluate AIDS

GLOBAL COHORTS COMPOSED OF  
**>1,000,000**  
PERSONS LIVING WITH HIV/AIDS

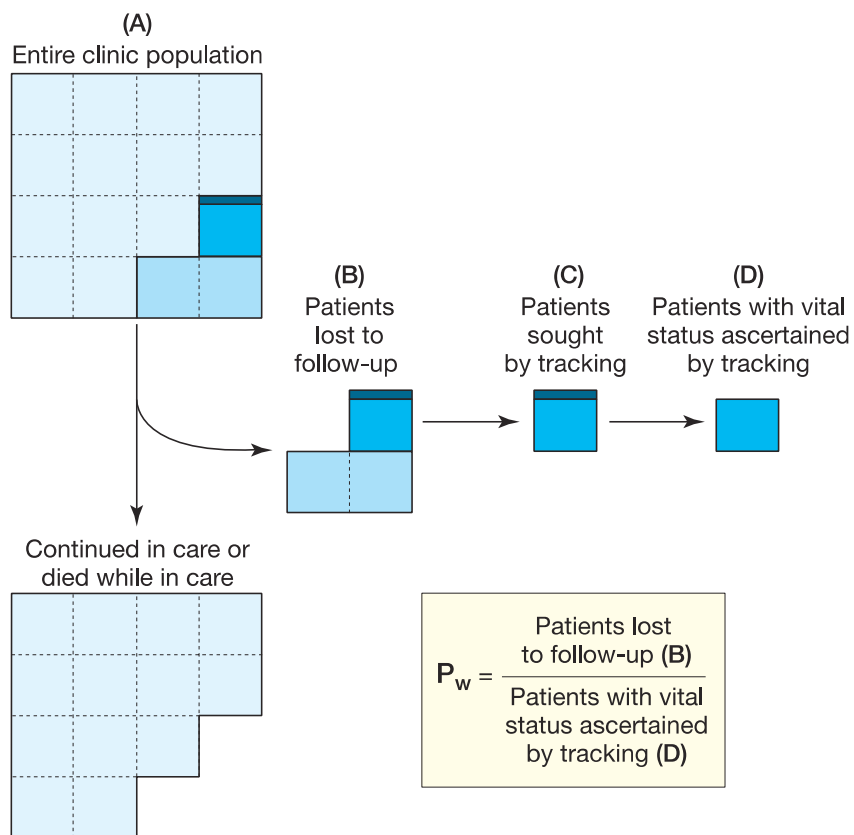


## Sampling-Based Approach to Determining Outcomes of Patients Lost to Follow-Up in Antiretroviral Therapy Scale-Up Programs in Africa

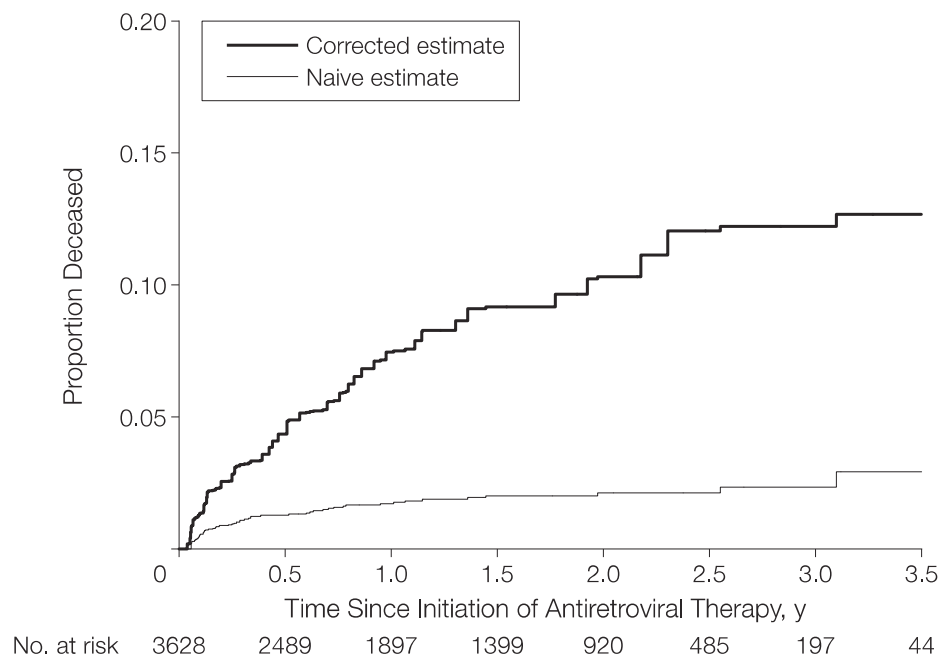
Elvin H. Geng; Nneka Emenyonu; Mwebesa Bosco Bwana; et al.

JAMA. 2008;300(5):506-507 (doi:10.1001/jama.300.5.506)

**Figure 1.** Derivation of Probability Weights



**Figure 2.** Naive and Corrected Mortality Estimates



# **PB N°3: LA CAUSALITÉ**

# Causalité

## Précautions

- Questions/hypothèses posées en amont
- Critères de causalité (Bradford Hill)
  - Critères internes à l'étude (force de l'association, relation temporelle, relation dose-réponse et spécificité de la relation)
  - Externes (consistance, cohérence, plausibilité biologique, analogie, expérimentation)
- Données -omiques
  - Connaissances des mécanismes biologiques
  - Validation des résultats par des nouvelles expérimentations spécifiques

# EN RECHERCHE CLINIQUE



ELSEVIER  
MASSON



Disponible en ligne sur

ScienceDirect  
[www.sciencedirect.com](http://www.sciencedirect.com)

Elsevier Masson France

EM|consulte  
[www.em-consulte.com](http://www.em-consulte.com)

Revue d'Épidémiologie et de Santé Publique 62 (2014) 1–4

---

---

Revue d'Épidémiologie  
et de Santé Publique

Epidemiology and Public Health

---

---

## Éditorial

### L'analyse des « Big Data » en recherche clinique

The analysis of “Big Data” in clinical research

R. Thiébaud <sup>abc,de,\*</sup>, B. Hejblum <sup>abc,de</sup>, L. Richert <sup>abc,de</sup>

<sup>a</sup>Inserm U897<sup>1</sup> épidémiologie et biostatistique, 33000 Bordeaux, France

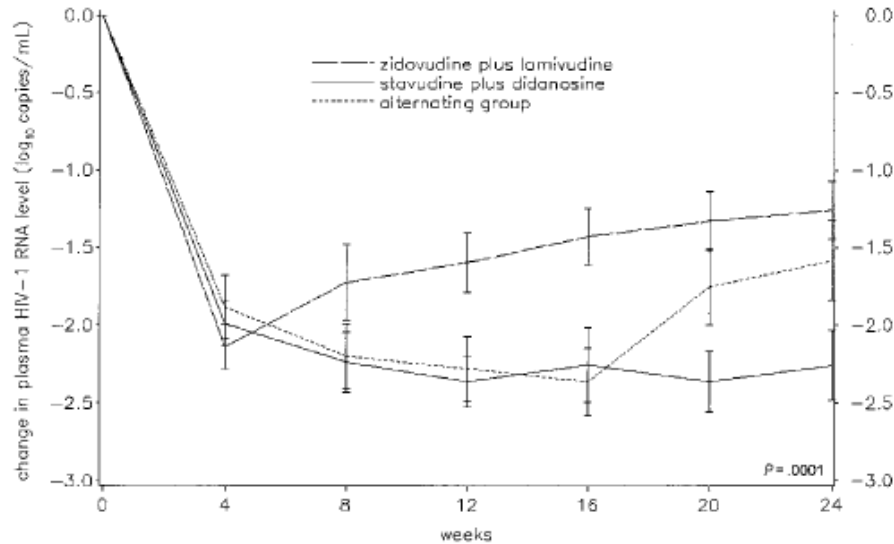
<sup>b</sup>Inria SISM<sup>2</sup>, 33000 Bordeaux, France

<sup>c</sup>Institut de santé publique d'épidémiologie et de développement (ISPED<sup>3</sup>), université Bordeaux, 33000 Bordeaux, France

<sup>d</sup>Vaccine Research Institute, 94010 Créteil, France

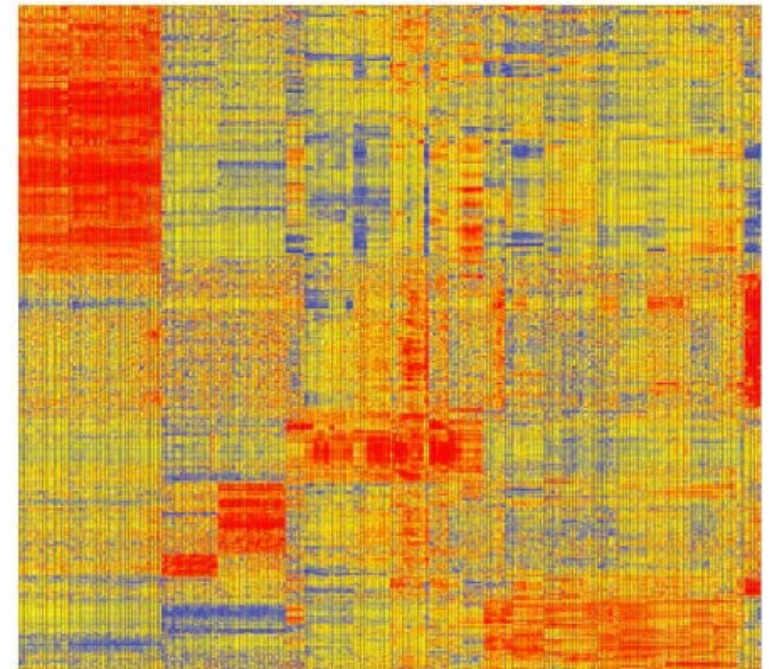
<sup>e</sup>Unité de soutien méthodologique à la recherche clinique et épidémiologique (USMR), CHU de Bordeaux, 33000 Bordeaux, France

# Exemple de deux essais cliniques



no. of patients

zidovudine-lamivudine	51	48	45	45	45	46	46
stavudine-didanosine	51	47	47	44	43	42	48
alternating group	49	47	44	42	39	47	46



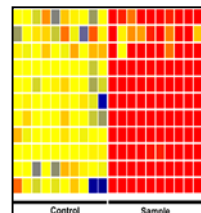
Essai	ALBI ANRS 070	DALIA-1
Nb de participants	151	19
Taille	67 Ko	200 Go



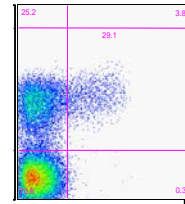
### Les (big) data

- 846 000 probes (18 temps x 47 000 sondes) 26 Mo
- 18 612 000 beads (22 billes/sonde) 6 Go

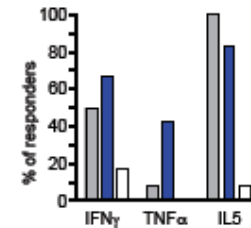
- 30 populations cellulaires 0.05 Mo
- 2160 anticorps (18 temps x 15 tubes x 8 anticorps) pour 2.6 Go



Gene profiling

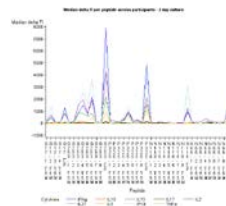


Cell responses

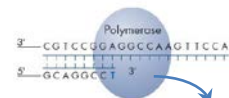
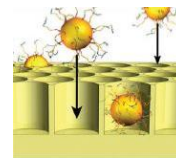


Cytokines

- 800 mesures/temps
- 0.35 Mo



Epitope mapping



PPi,  
H<sup>+</sup>,...

Viral changes/  
adaptation

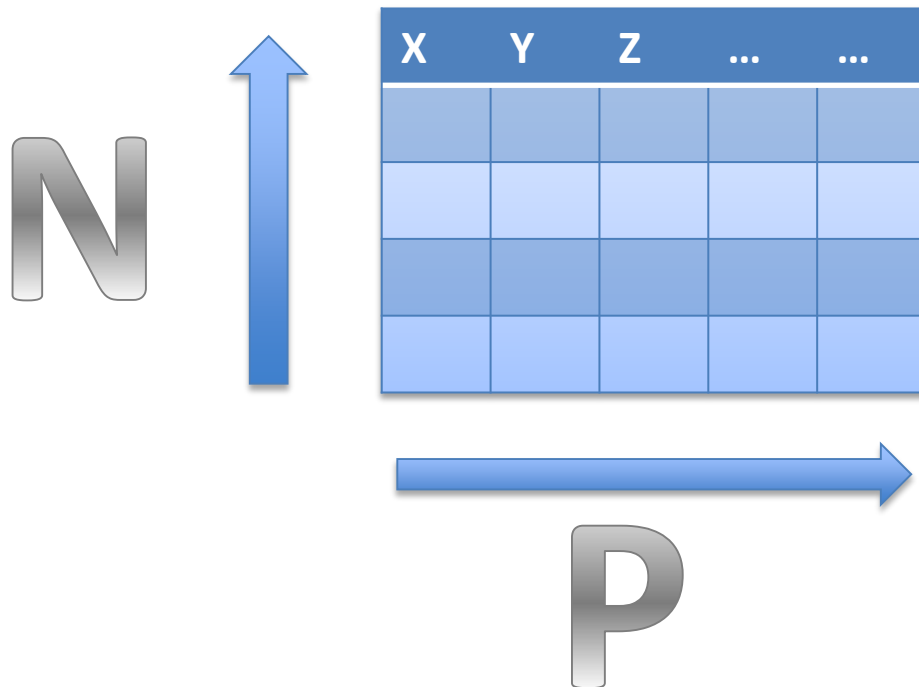
- 200 séquences
- 20 Mo

× 19

# Point de vue statistique

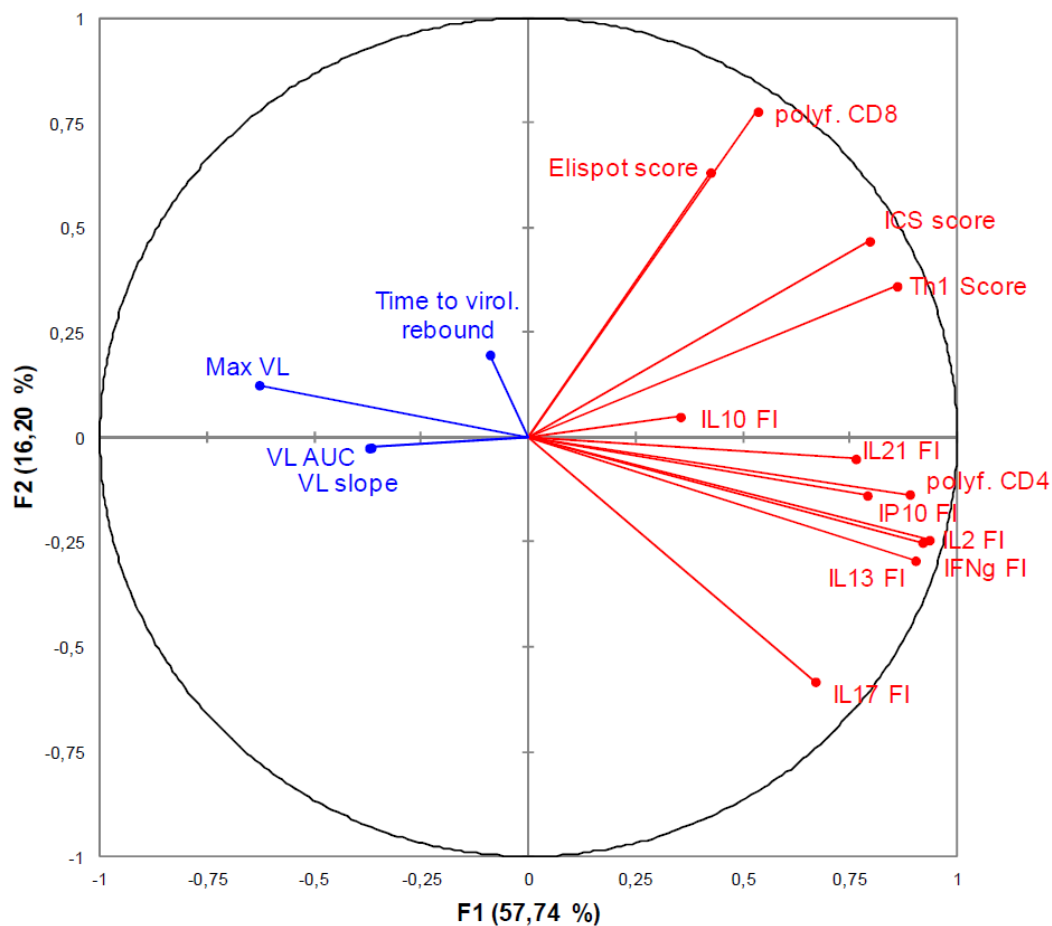
Données multidimensionnelles

$$N \ll P$$



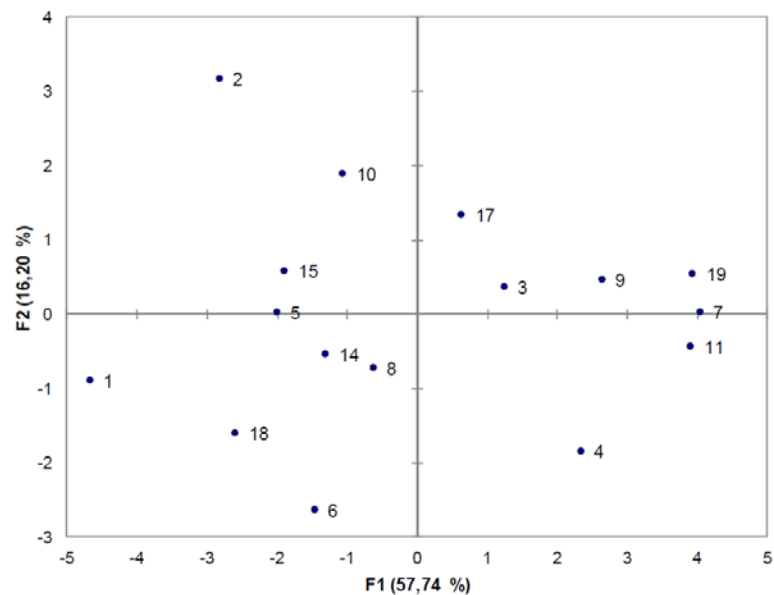
# Réponses immunologiques et virales

Variables (axes F1 and F2: 73,94 %)

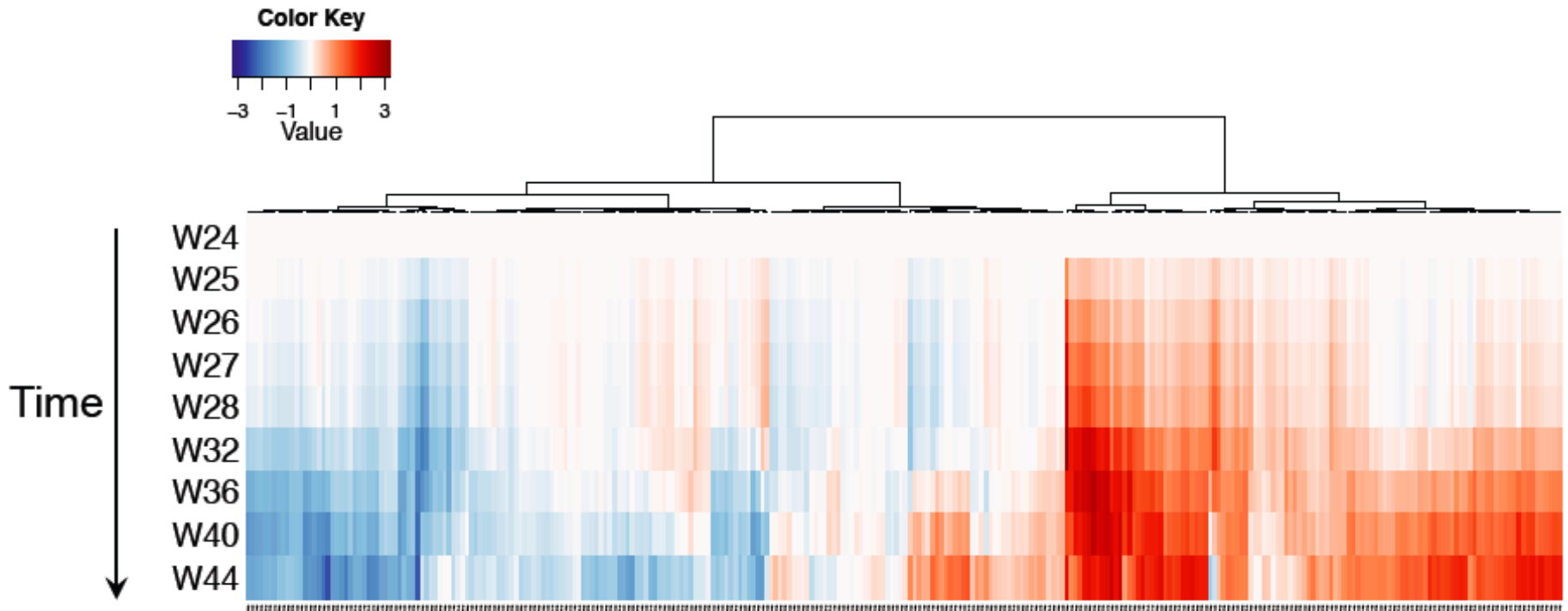


• Active variables    • Supplementary variables

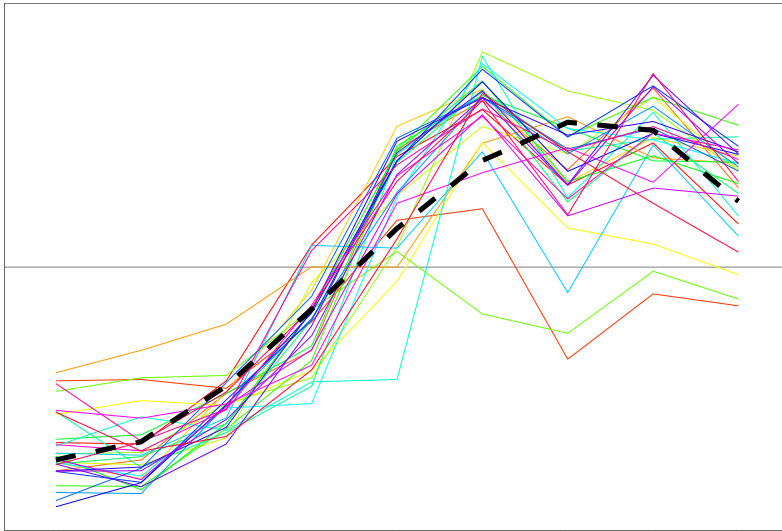
Individuals (axes F1 and F2: 73,94 %)



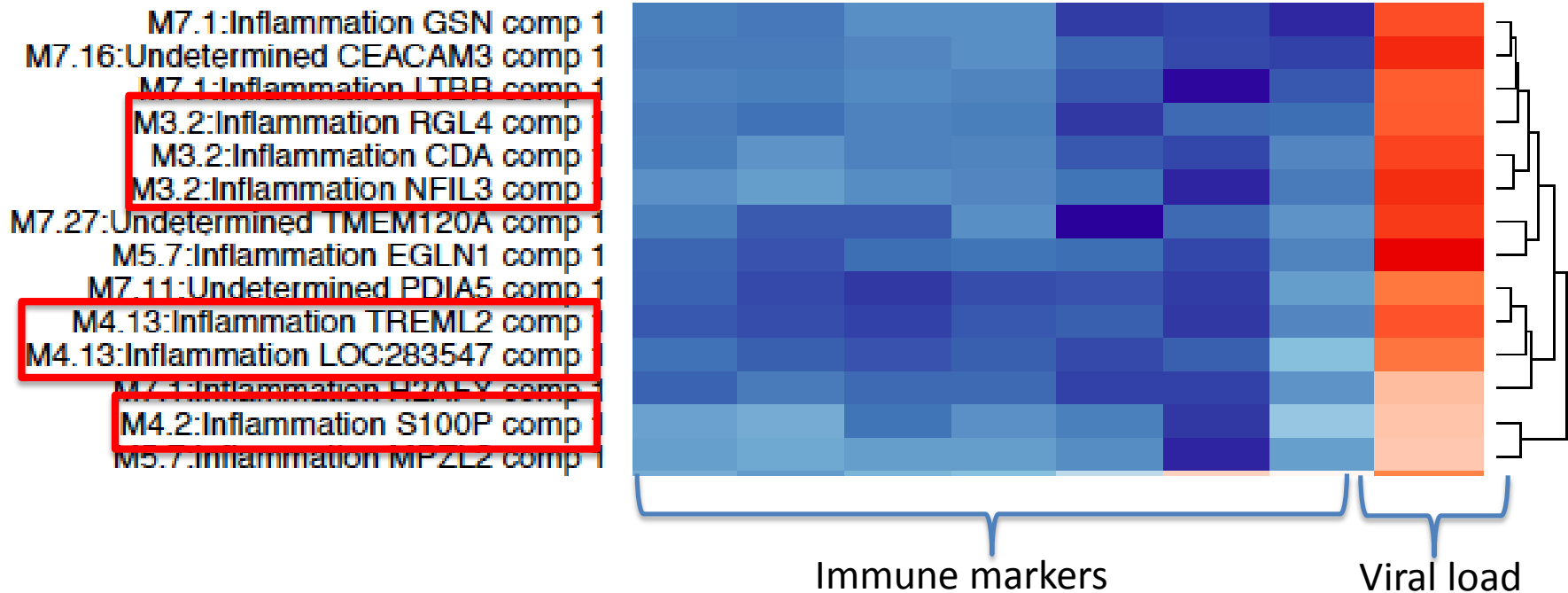
# Expression génique



# Expression génique: Time-course Gene Set Analysis



# Correlations entre groupes d'expression génique, réponses immunes (W16) et pic de charge virale



Association négative des groupes de gènes retrouvée aussi avec les réponses à la vaccination anti-pneumococciques (Obermoser et al., Immunity 2013)

# Conclusion: Big data en épidémiologie

## Deux situations

- Les données volumineuses issues
  - Des bases de données médico-administratives
  - Des collaborations multicohortes
  - Des consortium (GWAS)
    - ! facteurs non mesurés, puissance statistique
- Les données « riches » issues de nouvelles technologies
  - Données omiques
  - Données d'imagerie...
    - ! Fléau de la dimensionalité

## Conclusion: les 4 V



Volume

Variété

Vélocité

Validité

- Stockage, sécurité
- Débit
- Puissance de calcul

- **Comprendre les données**
- **Visualiser les données**
- **Analyser les données**<sup>31</sup>



Research article

Open Access

## Sparse canonical methods for biological data integration: application to a cross-platform study

Kim-Anh Lê Cao<sup>\*1,2</sup>, Pascal GP Martin<sup>3</sup>, Christèle Robert-Gr. Philippe Besse<sup>2</sup>

Journal of Biomedicine and Biotechnology • 2005:2 (2005) 147–154 • DOI: 10.1155/JBB.2005.147

RESEARCH ARTICLE



Disponible en ligne sur  
ScienceDirect  
www.sciencedirect.com

Elsevier Masson France  
EM|consulte  
www.em-consulte.com

Revue d'Épidémiologie  
et de Santé Publique  
Epidemiology and Public Health

Revue d'Épidémiologie et de Santé Publique 62 (2014) 1–4

Éditorial

### L'analyse des « Big Data » en recherche clinique

The analysis of "Big Data" in clinical research

R. Thiébaud<sup>abc,de\*</sup>, B. Hejblum<sup>abc,de</sup>, L. Richert<sup>abc,de</sup>

<sup>1</sup>Inserm U897<sup>1</sup> Épidémiologie et biostatistique, 33000 Bordeaux, France

<sup>2</sup>Inria SISIM<sup>2</sup>, 33000 Bordeaux, France

<sup>3</sup>Institut de santé publique et d'épidémiologie et de développement (ISPED<sup>3</sup>), université Bordeaux, 33000 Bordeaux, France

<sup>4</sup>Vaccine Research Institute, 94010 Créteil, France

<sup>5</sup>Unité de soutien méthodologique à la recherche clinique et épidémiologique (USMIS), CHU de Bordeaux, 33000 Bordeaux, France

BIOINFORMATICS

Vol. 00 no. 00 2005  
Pages 1–8

## Group and Sparse Group Partial Least Square Approaches Applied in Genomics Context

Benoît Liquet<sup>1,2,3,4,5\*</sup>, Pierre Lafaye de Micheaux<sup>6</sup>, Boris P. Hejblum<sup>2,3,4,5</sup> and Rodolphe Thiébaud<sup>2,3,4,5</sup>

<sup>1</sup>School of Mathematics and Physics, The University of Queensland, Brisbane 4066, Australia

<sup>2</sup>Inria, SISTM, Talence, France

<sup>3</sup>Inserm, U897, Bordeaux, France

<sup>4</sup>Bordeaux University, Bordeaux, France

<sup>5</sup>Vaccine Research Institute, Creteil, France

<sup>6</sup>Department of Mathematics and Statistics, Université de Montreal, Canada

## Classification and Selection of Biomarkers in Genomic Data Using LASSO

Debashis Ghosh<sup>1</sup> and Arul M. Chinnaiyan<sup>2</sup>

<sup>1</sup>Department of Biostatistics, University of Michigan, 1420 Washington Heights, Ann Arbor, MI 48109-2029, USA

<sup>2</sup>Departments of Pathology and Urology, University of Michigan, 1300 Catherine Road, Ann Arbor, MI 48109-1063, USA

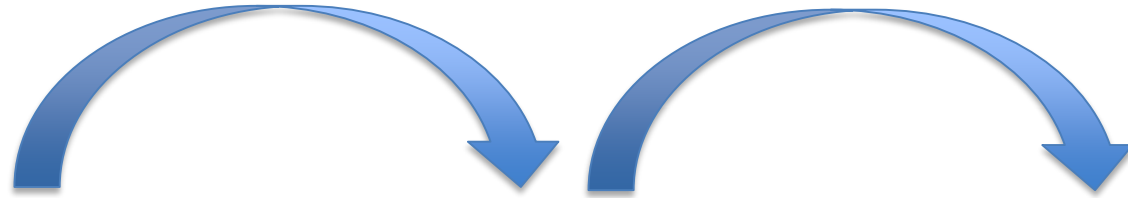
## Systems biology approaches to epidemiological studies of complex diseases

Hongzhe Li\*

WIREs Syst Biol Med 2013, 5:677–686.

# Qualité des données

- Garbage in, garbage out?



## Conclusion

- Big data in epidemiology = Big opportunity + Big challenge
- Besoin:
  - De multidisciplinarité
  - De méthodes d'analyse adaptées
  - Et du savoir faire épidémiologique et de l'intégrité scientifique

# Remerciements



Rodolphe Thiébaud



Linda Wittkop



Boris Hejblum



Agence autonome de l'Inserm

