

## The Exposome: Promises and Challenges of a New Concept

L'exposome : Promesses et défis d'un nouveau concept

**Rémy Slama** Collège de France & Inserm The relations between human health and the environment in the Anthropocene

Lecture #7 – 25 May 2022



1



### Lecture overview

- A. Motivation, aims & challenges
- *B. Layer 1:* exposome descriptive studies (biomonitoring / environmental justice challenge)
- C. 2-layer problem: exposome-health studies
- D. Multi-layer problem: Cross-omics analyses
- E. Perspectives: the first 20 years of exposome research

The French biomonitoring survey (seminar of Dr. Clémence Fillol)











## *Question 1:* What happens to health outcome c if E<sub>i</sub> varies? [classical aetiology]

The answer of « classical » risk factor epidemiology (assuming a randomized controlled experiment is not feasible, e.g. for ethical reasons):

- Recruit a cohort of "healthy" subjects
- Assess E (at inclusion and possibly repeatedly)
- Assess sociodemographic and behavioural factors also potentially influencing c
- Wait for outcome c to occur in a large enough number of subjects of the cohort
- Quantify the statistical relation between E and c adjusting for the "potential confounders"
- Repeat the study a certain number of times
- Perform a meta-analysis of all published studies and report to decision makers (if possible incorporating external mechanistic evidence and quantifying the overall level of evidence)





Figure 7 Survival of British doctors from age 35 years by smoking habits: lifelong nonsmokers—, cigarette smokers, smoking 1–14 cigarettes a day ..., smoking 15–24 cigarettes a day -..., 25 or move cigarettes a day ---- 1<sup>st</sup> (Reproduced with permission of the BMJ Publishing Group)

# Problems with the approach of classical risk factors epidemiology

- E<sub>i</sub> may be measured with error (measurement error)
- Confounders (and outcome) may be measured with error
- There may be confounding by other unmeasured exposures (confounding by coexposures)
- Effect-measure modifications ("interactions/synergy") not easily studied (unless strong a priori hypothesis)
- There may be hidden multiple testing (e.g., because actually E<sub>i</sub> was not an a priori choice)
- Random fluctuations
- Rather *low throughput* approach compared e.g. to toxicology (1 exposure every 1-2 years per PI in a given cohort?)
- The study is more likely to be published if it highlights a "significant" association, biasing the meta-analysis (publication bias)

#### Selling the salami by the slice... (i.e., studying each exposure separately in the same cohort)



Risk of selective reporting, publication bias, chance finding through multiple-testing...



Crab Uca pugnax (Wild, Cancer Epid Biomarkers Prev, 2005)





### Challenges of exposome studies

- 1. Exposure assessment: width vs. accuracy trade-off (can we assess more compounds without increasing exposure misclassification?)
- 2. Associating the exposome and health: multiple-testing issues (can we consider more exposures without increasing too much bias and random error, while maintaining the ambition of deciphering causal relations?)
- 3. Cross-exposure (mixtures/synergy) and cross-omics analyses



### Metrologic issue

Increasing the number of (exposure) factors considered should not be done at the cost of a decrease in the quality of their assessment. (cf. *curse of dimensionality* data science concept)

More exposures, better characterized

















# Linking the child *postnatal* exposome (125 exposures) with children lung function (FEV1 – Forced Expiratory Volume in 1s; 1033 children)

















# What happens when the dimension of a data set increases?

#### The 2 curses of dimension

High dimension data: when the number of variables p becomes large compared to n



## 1. The "real" curse of dimension: data sparsity



40

# 2. Increased false positive rate due to multiple testing

Remember that in exposome studies, the general aim is not a predictive one (can we predict Y given X<sub>1</sub>, ... X<sub>p</sub> ?) but rather a causal inference aim (which exposures among X<sub>1</sub>, ... X<sub>p</sub> causally influence Y?)



### An illustration of multiple testing issues



42

### Background: A very short history of multiple testing issues in social and health sciences (1)

- Quételet was probably the first to bring statistics from astronomy and meteorology into social and health sciences ("Research on population births, deaths, prisons, poor houses etc. in the kingdom of the Low Countries", 1827)
- Cournot-Quételet debate (1843)
   If one tests the sex ratio of all 86 French "départements" for a difference with
   the others, then the resulting p-values are meaningless
- In its early developments (1850-1980), epidemiology focused on "low dimension" problems with strong a priori hypotheses in which multiple testing was (at least apparently) not an issue

John Snow and the propagation of cholera, smoking and lung cancer...



Quételet (1797-1874)



Cournot (1801-1877)

# A very short history of multiple testing issues in social and health sciences (2) – The Genome era

- 1950s: First multiple comparison procedures
- 1979: Holm-Bonferroni method Family Wise Error Rates (FWER) Control the probability that NONE of the multiple observed scores are below a threshold  $\alpha$
- 1995: False Detection Rate (FDR) procedures
   Aims to make sure that the the *overall* rate of false positive signals remains below
   a threshold α. (Benjamini & Hochberg, J Roy Stat Soc B, 1995) for independent tests, Benjamini &
   Yekutieli (Ann Stat, 2001) under arbitrary dependence
- 1998: 1<sup>st</sup> commercial Affymetrix array (1494 single nucleotide polymorphisms SNPs)
- 2016: Genomic arrays feature 1.8 million genetic markers
- These multiple comparison procedures are now widely used in the genomic (and "epigenomic") literature











# Are there efficient approaches to detect *interactions*\* between exposures?

Sensitivity to detect order-2 interaction terms



False positive rate for interaction terms







#### Issues related to power and sample size

54

### The curse of measurement error









f (sample size (or number of cases), exposure distribution, measurement error...)

Related to the within-subject variability of the compound





How many (pools) of urine samples are enough to accurately assess exposure to non persistent compounds?

Estimation of the number of urine pools needed to limit exposure misclassification of phthalate metabolites, phenols, OP pesticide metabolites, and cotinine in pregnant women and children based on the intraclass correlation coefficients

Biomarkers	Pregnant women	Children	
	N of pools of 20 urines needed for between-trimester $ICC \ge 0.80$	N of pools of two urines needed for between-week ICC $\ge 0.80$	N of pools of 15 urines needed for between-season ICC $\geq 0.80$
Phthalate me	tabolites		
MEP	3	3	4
MiBP	3	3	3
MnBP	2	3	4
MBzP	3	3	3
MEHP	3	3	2
MEHHP	4	3	3
MEOHP	4	3	3
MECPP	3	3	3
oh-MiNP	4	3	4
oxo-MiNP	4	4	4
Phenols			
MEPA	4	3	4
ETPA	3	3	4
PRPA	4	3	4
BUPA	3	3	4
BPA	4	3	5
OXBE	3	1	3
TCS	3	1	5
OP pesticide	metabolites		
DMP	4	4	5
DMTP	4	4	4
DMDTP	6	4	4
DEP	3	4	4
DETP	4	4	4
DEDTP	4	4	4
Cotinine	2	-	3











### Example: HELIX project has characterized multiple layers 'Omics signatures in 1200 children



#### In 1200 children (6-10 years)

Platform	# features		
Urine metabolome	44		
Serum metabolome	177		
Proteome	36		
miRNAs	359		
Transcriptome	35,841		
Methylome	386,518		
Exposome	Ca. 120		
Health	Ca. 10		

(Gallego-Paüls, BMC Medicine 2021)

Exposome studies incorporating multiple layers of intermediary biological ('omics) data: *Motivation* 



(technically realistic?)

#### Solving more easily the 2-layer problem of exposome-health relations by borrowing information from another intermediary biological layer





# Exposome studies incorporating multiple layers of intermediary biological ('omics) data: *Motivation*

- Increase the level of evidence for the exposome (or a specific set of exposures) influencing disease risk, through identification of toxicologicallyplausible biological mediators
- 2. Adopt a "systems biology" approach, i.e. decipher the causal relations among the multiple biological layers (high dimension multilayer causal inference)
- 3. Solve more easily the 2-layer problem of exposome-health relations by borrowing information from other intermediary biological layers
  - Dimension reduction
  - Discard reverse causality
  - Identify exposure biomarkers or effect biomarkers
- 4. Disease **risk prediction** (without necessarily pretention to causal inference)

Causal inference in a multi-layer (≥3) setting (very ambitious)

Causal inference in a 2-layer setting (ambitious but realistic)

> Risk prediction (technically realistic?)

### High-dimension mediation analysis

- Mediation analysis is conceptually well framed with a mediator of dimension 1 (VanderWeele, *Oxford Univ Press*, 2015). The theoretical framework in particular implies that the causal model is known a priori.
- Such an assumption is not very realistic if the mediator has a high dimension and is treated as a large set of potential mediators
  - Remember that *statistically*, models generally do not allow to infer the direction of any causal effect between A and B (i.e., A->B and B->A are not distinguishable by pure statistical tools)
  - And that biologically there may be complex causal relations within a biological layer (e.g., the methylome)
- Consequently, although the literature is full of examples of high dimension analysis relying e.g., on methylome or metabolome data, rigorously identifying the "causal" mediators or the share of the effect of E on Y mediated by an intermediary biological layer generally remains a challenge.



VanderWeele, Epid Meth, 2014; Blum, EHP, 2020)

# Exposome studies incorporating multiple layers of intermediary biological ('omics) data: *Motivation*

- 1. Increase the level of evidence for the exposome (or a specific set of exposures) influencing disease risk, through identification of toxicologically-plausible biological mediators
- Adopt a "systems biology" approach, i.e. decipher the causal relations among the multiple biological layers (high dimension multilayer causal inference)
- 3. Solve more easily the 2-layer problem of exposome-health relations by borrowing information from other intermediary biological layers
  - Dimension reduction
  - Discard reverse causality
  - Identify exposure biomarkers or effect biomarkers
- Disease risk prediction (without necessarily claim to causal inference)

Causal inference in a multi-layer (≥3) setting (very ambitious)

Causal inference in a 2-layer setting (ambitious but realistic)

> Risk prediction (technically realistic?)

### Disease risk prediction using 'omics data

- Disease risk prediction (with claim to identifying causal health predictors) is not particularly
  relevant for exposome research, in which one generally aims to identify actionable environmental
  disease drivers likely to allow public health improvement
  - *Exception*: identification of predictors of exposures from omics (e.g., methylome) signals (Guida, *Mol Hum Gen*, 2015)
- Statistical learning tools perform generally well when it comes to prediction (as opposed to causal inference).
- However, they tend to do so when the number of "training samples" is large, which is typically not the case currently for 'omics studies in the health field, which are generally conducted on a low number of subjects (n ca. 10<sup>3</sup>-10<sup>4</sup>, for a number of features generally in the 10<sup>4</sup>-10<sup>8</sup> range)
- Internal validation of models provides over-optimistic estimations of the classification accuracy (i.e., overfitting). Leave one out cross-validation seems particularly prone to such overconfidence in the predictive ability (Rodriguez-Perez, *Anal Bioanal Chem*, 2018)

#### Omics markers are also highly variable

- Intra-, inter-individual and cohort variability of multiomics profiles **measured 6 months apart** in 156 children

- DNA methylation most stable; expression, least stable
- Strong heterogeneity between features







# The historian's corner: A parallel between genome and exposome research



79



07/09/2022

## Cohorts of connected study participants

Enki Bilal







### Growing in number or drowning by numbers?





Not safe to increase the number of exposures considered... -if you cannot simultaneously improve the quality of their assessment (which can be done by increasing the number of biospecimens collected per subject; see Perrier, *Epidemiology*, 2016) -if you cannot simultaneously increase sample size

86



#### Some possible questions

- 1) What happens to health outcome c if E<sub>i</sub> varies? (e.g., : does smoking cause lung cancer?) [classical, single exposure, aetiology]
- What are all the external risk factors that influence c? [exposome-health study/exposome-wide study]
- What happens to health if E<sub>i</sub> varies? (e.g., : what are the sanitary consequences of smoking?) [outcome-wide aetiology]
- 4) What intermediary variables may explain an effect of E<sub>i</sub> on c? [Mechanistic research/mediation analysis]
- 5) Which risk factor has the largest *impact* (i.e. attributable number of disease cases) on c today? [Environmental burden of disease] On health overall?
- 6) Are specific sociodemographic subgroups disproportionately exposed to harmful exposures? [Environmental justice problem]

Still a challenge (power, FDP) Feasible (but not an exposome study) Conceptual issues Complex

Within-reach

# The *Exposome/environmental health* ambitious road

		Exp	osome			
Quantification	Link with socio- territorial characteristics « Environmental justice »		Health e In vivo toxic Cohorts	effects ology (AOPs),	Act on Individu Governar framewor	the exposome al interventions ace / regulatory k
Exposome-ready » cohor	ts	Mechan Cross-omi In vitro and cohorts	hisms of action cs studies in vivo toxicology,		th impac nmental dis	t sease burden
	M	ethodo	logical issues	5		

### As a conclusion

- The exposome is a challenging but fruitful concept for environmental health research
- It allows to explicitly tackle essential challenges sometimes hidden in single exposure environmental epidemiology
  - Publication bias, differential measurement error across exposures...
- Exposome studies shouldn't consist in studies with a number of subjects is similar (or lower, as is often the case) to single exposure studies, with more exposures assessed.
  - Design issues, to some extent already faced when genetic epidemiology went genome-wide
- The (probably less challenging) outcome-wide approach (VanderWeele, Epidemiology, 2016) is also worth considering when it comes to improving the throughput of human studies.
- Epidemiology and toxicology need each other to walk along the way of the exposome
  E.g., using toxicology to a priori reduce the dimension of epidemiological exposome studies