

Understanding complex data

Pietro Ferrari, PhD

Head, Nutritional Methodology and Biostatistics Group IARC, Lyon, France

1 / 38

Outline

- Metabolomics as an emerging resources
- The PC-PR2 technique to investigate drivers of variability in large dimension data
- Biomarker discovery strategy
- The meeting-in-the-middle (MITM) principle

- Applications from the EPIC study



Metabolomics

• Comprehensive and quantitative analysis of wide arrays of metabolites in biological samples

イロト イポト イヨト イヨト 二日

- Progressively available in prospective epidemiological investigations



Metabolomics

- Comprehensive and quantitative analysis of wide arrays of metabolites in biological samples
- Progressively available in prospective epidemiological investigations
- Powerful tool for potential identification of causal pathways in disease development
- Need for statistical methodologies



The EPIC Study

- Prospective cohort with 500,000 participants from 23 centres
 - Dietary and lifestyle exposures assessed at baseline
 - Biological samples collected at baseline from 80% disease-free participants





EPIC study on HCC

- Nested case-control study on hepatocellular carcinoma (HCC)
- 147 cases and 147 matched controls
- 132 blood metabolites acquired in cancer-free individuals with Biocrates AbsoluteIDQ-p180 Kit using UPLC coupled to mass spectrometer





Explore variation in -omics data

• Before doing anything else, spend time exploring variability in metabolomics data



Biocrates metabolites



World Health Organization

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Explore variation in -omics data

- Before doing anything else, spend time exploring variability in metabolomics data
- What drives systematic variation in -omics?



Explore variation in -omics data

- Before doing anything else, spend time exploring variability in metabolomics data
- What drives systematic variation in -omics?
- Subjects general characteristics (e.g. age, BMI, smoking status)
- Variables describing samples' technical treatment (e.g., serum clot contact time, length of storage, fasting status)



The structure of -omics data



Standard scenario





New generation scenario







11/38





- Run principle component analysis (PCA)
- Retain the first $q~(q \ll p)$ components



Displaying PCA structure



イロン 不同と 不同と 不同と



A new method: PC-PR2

- Only two components at the time are displayed
- Only one explanatory factor can be visualised

ヘロン 人間と 人間と 人間と

14/38

- Important inter-correlations between explanatory variables NOT accounted for



A new method: PC-PR2

- Only two components at the time are displayed
- Only one explanatory factor can be visualised

・ロン ・回 と ・ 回 と ・ 回 と

14/38

- Important inter-correlations between explanatory variables NOT accounted for
- Principal Component Partial R-square (PC-PR2) method developed



Analytical steps of PC-PR2

(a) Perform PCA on the -omics data

15/38



Analytical steps of PC-PR2

- (a) Perform PCA on the -omics data
- (b) Retain **q** components, explaining -omics variability above a given threshold, ie. 80%

イロト イポト イヨト

15/38



Analytical steps of PC-PR2

- (a) Perform PCA on the -omics data
- (b) Retain **q** components, explaining -omics variability above a given threshold, ie. 80%
- (c) Fit **q** linear regression models, where each component (dependent variables) is explained in terms of **k** explanatory variables



Step (c) of PC-PR2

• For each component q with $(q = 1, ..., q_{tot})$

$$\mathsf{PC}_{qi} = \beta_0 + \beta_1 X_{1i} + \ldots + \beta_k X_{ki} + \epsilon_{qi}$$

with (i=1,..., n) and $\epsilon_{qi} \sim \mathsf{N}(0, \sigma_{\epsilon q}^2)$,

イロト イヨト イヨト イヨト 三日

16/38



Analytical steps of PC-PR2 (ii)

- (c) Fit **q** linear regression models, where each component (dependent variables) is explained in terms of **k** explanatory variables
- (d) Determine the $R^2_{partial(q,k)}$ for each explanatory variable, in each model q

イロン 不良 とくほど 不良 とうほ



Step (d) of PC-PR2

The partial R-square (R²_{partial}) is a statistics quantifying the amount of variability of the dependent variable that each explanatory variable X_k contributes to explain, conditional on the other (k-1) variables in the model

$$R^2_{partial(\mathbf{X_1})q}$$



Step (d) of PC-PR2

The partial R-square (R²_{partial}) is a statistics quantifying the amount of variability of the dependent variable that each explanatory variable X_k contributes to explain, conditional on the other (k-1) variables in the model

$$R^2_{partial(\mathbf{X_1})q}, \ R^2_{partial(\mathbf{X_2})q}$$



Step (d) of PC-PR2

The partial R-square (R²_{partial}) is a statistics quantifying the amount of variability of the dependent variable that each explanatory variable X_k contributes to explain, conditional on the other (k-1) variables in the model

$$R^2_{partial(\mathbf{X_1})q}, R^2_{partial(\mathbf{X_2})q}, \ldots, R^2_{partial(\mathbf{X_k})q}$$



Analytical steps of PC-PR2 (ii)

- (c) Fit **q** linear regression models, where each component (dependent variables) is explained in terms of **k** explanatory variables
- (d) Determine the $R^2_{partial(q,k)}$ for each explanatory variable, in each model q
- (e) For each **k** variable, determine an overall $R_{partial(k)}^2$ as a weighted average, using eigenvalues as weights



Step (e) of PC-PR2

• For each variable **k**, over the **q** PCs:

イロト イヨト イヨト イヨト 三日

22/38



Step (e) of PC-PR2

• For each variable **k**, over the **q** PCs:

イロト イヨト イヨト イヨト 三日

22/38

$$R^2_{ ext{partial}(old X_1)} = \sum R^2_{ ext{partial}(old X_1)q} w_q$$



Step (e) of PC-PR2

• For each variable **k**, over the **q** PCs:

イロト イヨト イヨト イヨト 三日

22/38

$$R_{partial(\mathbf{X}_{1})}^{2} = \sum R_{partial(\mathbf{X}_{1})q}^{2} w_{q}$$
$$R_{partial(\mathbf{X}_{2})}^{2} = \sum R_{partial(\mathbf{X}_{2})q}^{2} w_{q}$$



Step (e) of PC-PR2

• For each variable **k**, over the **q** PCs:

$$R_{partial(\mathbf{X}_{1})}^{2} = \sum R_{partial(\mathbf{X}_{1})q}^{2} w_{q}$$
$$R_{partial(\mathbf{X}_{2})}^{2} = \sum R_{partial(\mathbf{X}_{2})q}^{2} w_{q}$$
$$\dots$$

$$R^2_{\it partial(f X_K)} = \sum R^2_{\it partial(f X_K)q} w_q$$

イロト イヨト イヨト イヨト 三日

22/38



Step (e) of PC-PR2 (ii)

• For each variable **k**, over the **q** PCs:

$$R_{partial(\mathbf{X}_{\mathbf{K}})}^{2} = \sum R_{partial(\mathbf{X}_{\mathbf{K}})q}^{2} w_{q}$$

- where
$$\textit{w}_{q} = \lambda_{q} / \sum \lambda_{q}$$
, for $(\textit{q} = 1, \ldots, \textit{q}_{tot})$

with λ_q = eigenvalue that expresses the amount of variability captured by PC_q



The PC-PR2 analysis



Biomarker discovery



International Agency for Research on Cancer

(ロ) (同) (E) (E) (E)



Biomarker discovery



(ロ) (同) (E) (E) (E)



Biomarker discovery



25/38



Diet and polyphenols, n=475

$Exposure^1$	Selected PP ²	\hat{r}_{adj}	AUC ³ (95% CI)
Coffee	Single PP (caffeic acid)	0.42	86% (78%, 94%)
	PP by LASSO (p=11)	0.51	89% (83%, 95%)
Red wine	Single PP (Gallic acid EE)	0.66	89% (84%, 95%)
		0.66	

¹ 24-hour dietary recall measurements; ²Urinary poliphenols metabolites measured by UPLC-ESI-MS/MS in 24-hour urine; ³Estimated by cross-validation in test and training sets.

(Noh *et al.*, submitted to JN)



Biomarker discovery (ii)

$Exposure^1$	Selected PP ²	\hat{r}_{adj}	AUC ³ (95% CI)
Coffee	Single PP (caffeic acid)	0.42	86% (78%, 94%)
	PP by LASSO (p=11)	0.51	89% (83%, 95%)
Red wine	Single PP (Gallic acid EE)	0.66	89% (84%, 95%)
	PP by LASSO (p=2)	0.66	89% (84%, 95%)

¹ 24-hour dietary recall measurements; ²Urinary poliphenols metabolites measured by UPLC-ESI-MS/MS in 24-hour urine; ³Estimated by cross-validation in test and training sets.

(Noh *et al.*, submitted to JN)





International Agency for Research on Cancer



(Vineis & Perera, CEBP, 2007)

28 / 38

Meeting-in-the-middle (MITM)



(Assi et al., Mutagenesis, 2015)

- T 🗎

29/38





・ロン ・回 と ・ ヨ と ・ ヨ と

30/38



Analytical strategy of MITM **Metabolomics** 2. Etiology

BMI HCC

・ロン ・回 と ・ ヨ と ・ ヨ と



Analytical strategy of MITM



・ロン ・四マ ・ビア・ ・ロン



1. BMI-driven PLS factor

Metabolites	Loadings	
Glutamine	-0.19	
Glutamate	0.23	
Tyrosine	0.24	
Lyso PC a C17:0	-0.22	
Lyso PC a C18:2	-0.23	
PC ae C36:2	-0.20	
Liver function score	0.19	

International Agency for Research on Cancer

(Assi *et al.*, submitted to PLOS Medicine). $\frac{33/38}{33}$

2. BMI and HCC

Exposure	OR ¹ (95%CI)	p_{value}
BMI	1.23 (0.93, 1.62)	0.149
Metabolites, PLS score	4.04 (2.22, 7.36)	4.8E-07
% Mediated ²	100	

¹ Expressing HCC relative risk estimate of 1-SD increase in the PLS score; ² Estimated as In(NIE)/[In(NIE)+In(NDE)], with NIE=natural indirect effect, NDE=natural direct effect.

International Agency for Research on ConArssi et al., submitted to PLOS Medicine)



2. BMI and HCC

Exposure	OR ¹ (95%CI)	p_{value}
BMI	1.23 (0.93, 1.62)	0.149
Metabolites, PLS score	4.04 (2.22, 7.36)	4.8E-07
% Mediated ²	100	

¹ Expressing HCC relative risk estimate of 1-SD increase in the PLS score; ² Estimated as In(NIE)/[In(NIE)+In(NDE)], with NIE=natural indirect effect, NDE=natural direct effect.

International Agency for Research on ConAssi et al., submitted to PLOS Medicine)



Limitations

- Sample size, accuracy, reliability
- Mediation analysis offers a framework to link inter-correlated factors, but ...

・ロト ・回ト ・ヨト ・ヨト

36 / 38



Limitations

- Sample size, accuracy, reliability
- Mediation analysis offers a framework to link inter-correlated factors, but ...

ヘロン 人間と 人間と 人間と

36 / 38

 It involves massive use of underlying assumptions, i.e. chronological sequence, confounding structure



Concluding remarks

- Strategies that tackle the complexity of dietary (and lifestyle) exposures are to be commended
- Vastly unexplored potential of -omics data
- But again, clearly biology is way more complex than statistical modeling
- Key to create multi-disciplinary settings, with a common language

(日) (部) (注) (注) (言)

37 / 38



Acknowledgments

- Anne Fages, Hwayoung Noh, Nada Assi (NMB Group), Vivian Viallon (Lyon I University)
- Pekka Keski-Rahkonen, Augustin Scalbert, Marc Gunter, Mazda Jenab, (NME Section)

- Paolo Vineis and all EPIC PIs

