

1 **A Statistical framework to model the meeting-in-the-middle principle using**
2 **metabolomic data: application to hepatocellular carcinoma in the EPIC study.**

3 *Nada Assi¹, Anne Fages^{2,a}, Paolo Vineis³, Marc Chadeau-Hyam³, Magdalena Stepien¹, Talita*
4 *Duarte-Salles¹, Graham Byrnes¹, Houda Boumazza², Sven Knüppel⁴, Tilman Kühn⁵, Domenico*
5 *Palli⁶, Christina Bamia⁷, Hendriek Boshuizen⁸, Catalina Bonet⁹, Kim Overvad¹⁰, Mattias*
6 *Johansson^{1,11}, Ruth Travis¹², Marc J Gunter³, Eiliv Lund¹³, Laure Dossus¹⁴⁻¹⁵, Bénédicte*
7 *Elena-Herrmann², Mazda Jenab¹, Vivian Viallon^{16-18†}, Pietro Ferrari^{1†*}*

8
9 ¹ International Agency for Research in Cancer (IARC-WHO), 150 Cours Albert Thomas,
10 69372 Lyon CEDEX 08, France

11 ² Centre de RMN à Très Hauts Champs, Institut des Sciences Analytiques (CNRS/ENS
12 Lyon/UCB Lyon 1), Université de Lyon, 69100 Villeurbanne, France

13 ³ Department of Epidemiology and Biostatistics, MRC-HPA Centre for Environment and
14 Health, School of Public Health, Imperial College London, Norfolk Place, London, W2 1PG,
15 United Kingdom

16 ⁴ Department of Epidemiology, German Institute of Human Nutrition, Potsdam-Rehbruecke,
17 14558 Nuthetal, Germany

18 ⁵ Division of Cancer Epidemiology, German Cancer Research Center (DKFZ), Heidelberg,
19 Germany

20 ⁶ Molecular and Nutritional Epidemiology Unit, Cancer Research and Prevention Institute –
21 ISPO, Florence- Italy

22 ⁷ WHO Collaborating Center for Food and Nutrition Policies, Department of Hygiene,
23 Epidemiology and Medical Statistics, University of Athens Medical School, Athens, Greece

24 ⁸ National Institute for Public Health and the Environment (RIVM), Bilthoven, The
25 Netherlands

26 ⁹ Unit of Nutrition and Cancer, Cancer Epidemiology Research Program, Institut Català
27 d'Oncologia, L'Hospitalet de Llobregat, Spain

28 ¹⁰ The Department of Epidemiology, School of Public Health, Aarhus University, Aarhus,
29 Denmark

30 ¹¹ The Department for Biobank Research, Umeå University, Sweden

31 ¹² Cancer Epidemiology Unit, Nuffield Department of Population Health University of
32 Oxford, Oxford, UK

33 ¹³ The Institute of Community Medicine, University of Tromsø, Tromsø, Norway

34 ¹⁴ Inserm, Centre for research in Epidemiology and Population Health (CESP), U1018,
35 Lifestyle, genes and health: integrative trans-generational epidemiology team, Villejuif,
36 France

37 ¹⁵ Université Paris Sud, Villejuif, France

38 ¹⁶ Université de Lyon, F-69622, Lyon, France

39 ¹⁷ Université Lyon 1, UMRESTTE, F-69373 Lyon, France

40 ¹⁸ IFSTTAR, UMRESTTE, F-69675 Bron, France

41 ^a Present address: Chemical Physics Department, Weizmann Institute of Science, Rehovot,
42 Israel

43 [†] Joint last authors.

44

45 ^{*}Corresponding author: Pietro Ferrari, International Agency for Research on Cancer, 150
46 Cours Albert Thomas, 69372 Lyon cedex 08, France. Tel: +33 472 73 8031; Fax: +33 472 73
47 8361. E-mail: ferrari@iarc.fr

48 **Abstract**

49 Metabolomics is a potentially powerful tool for identification of biomarkers associated with
50 lifestyle exposures and risk of various diseases. This is the rationale of the “meeting-in-the-
51 middle” concept, for which an analytical framework was developed in this study. In a nested
52 case-control study on hepatocellular carcinoma (HCC) within the European Prospective
53 Investigation into Cancer and nutrition (EPIC), serum ¹H NMR spectra (800 MHz) were
54 acquired for 114 cases and 222 matched controls. Through Partial Least Square (PLS)
55 analysis, 21 lifestyle variables (the “predictors”, including information on diet, anthropometry
56 and clinical characteristics) were linked to a set of 285 metabolic variables (the “responses”).
57 The three resulting scores were related to HCC risk by means of conditional logistic
58 regressions. The first PLS factor was not associated with HCC risk. The second PLS
59 metabolomic factor was positively associated with tyrosine and glucose, and was related to a
60 significantly increased HCC risk with OR= 1.11 (95%CI: 1.02, 1.22, p=0.02) for a 1-SD
61 change in the responses score, and a similar association was found for the corresponding
62 lifestyle component of the factor. The third PLS lifestyle factor was associated with lifetime
63 alcohol consumption, hepatitis and smoking, and had negative loadings on vegetables intake.
64 Its metabolomic counterpart displayed positive loadings on ethanol, glutamate and
65 phenylalanine. These factors were positively and statistically significantly associated with
66 HCC risk, with 1.37 (1.05, 1.79, p=0.02) and 1.22 (1.04, 1.44, p=0.01), respectively. Evidence
67 of mediation was found in both the second and third PLS factors, where the metabolomic
68 signals mediated the relation between the lifestyle component and HCC outcome. This study
69 devised a way to bridge lifestyle variables to HCC risk through NMR metabolomics data.
70 This implementation of the “meeting-in-the-middle” approach finds natural applications in
71 settings characterized by high-dimensional data, increasingly frequent in the –omics
72 generation.

- 73 **Keyword:** partial least square, lifestyle factors, metabolomics, hepatocellular carcinoma,
74 meeting-in-the-middle, molecular epidemiology.

75 **Introduction**

76 Metabolomic profiles from blood and other biological samples collected from large-
77 scale epidemiologic studies are increasingly being investigated [1], following recent
78 developments in nuclear magnetic resonance (NMR) and mass spectrometry (MS) enabling
79 the assessment of metabolic profiles for large numbers of individuals. As a result,
80 metabolomic data is gradually playing a key part in clinical and observational studies; and
81 new statistical methodologies [2] are increasingly being sought to explore insights into
82 pathological processes that metabolomics may provide in order to better understand
83 determinants of disease development. These approaches explore a variety of etiological
84 hypotheses; however they usually focus on one aspect at a time, combining metabolomics
85 with either epidemiologic/phenotypic data on lifestyle exposures [3] or with disease outcomes
86 [4,5]. The main aim of this work is to jointly use all aspects that are potentially informative to
87 apprehend the contrivances of disease development.

88 Metabolomic data offers the opportunity to identify signatures and biomarkers
89 associated with environmental exposures and the risk of a disease. Prospective studies are
90 conceptually suitable for this purpose, since they rely on biological samples collected before
91 disease onset, and are thus marginally influenced by metabolic changes due to processes of
92 disease development. In this scenario, the “meeting-in-the-middle” (MITM) approach [6] has
93 been conceived as a research strategy to identify biomarkers that are related to specific
94 exposures and that are, at the same time, predictive of disease outcome. Finding this overlap
95 between exposure and disease of “intermediate” biomarkers can potentially disclose useful
96 information on the exposure-to-disease pathway, and may serve as an objective risk exposure
97 measure, ultimately allowing the identification of a targeted prevention scheme. The MITM
98 was previously implemented as a proof of concept in a case-control study nested within a
99 cohort of healthy individuals [7], where a list of putative intermediate ¹H NMR biomarkers

100 linking exposure to dietary compounds, mainly micro- and macronutrients, and disease
101 outcomes (colon and breast cancer) were investigated.

102 In this study we extend previous attempts to model the MITM by fully integrating
103 metabolomics, lifestyle and disease risk in a single analytical framework. A strategy was
104 developed to simultaneously investigate a broad range of metabolites and lifestyle variables
105 with a partial least square (PLS) regression model [8]. The resulting scores were related to the
106 risk of hepatocellular carcinoma (HCC), in a case-control study nested within the European
107 Prospective Investigation into Cancer and nutrition (EPIC). HCC is the most frequent primary
108 form of cancer affecting the liver, an organ that plays a critical role in many metabolic
109 pathways [9]. HCC is a disease with multifactorial origins embracing lifestyle and dietary
110 exposures whose intersection may reveal metabolomic signals [10] relevant to cancer onset.
111 The system of relationships between metabolomic profiles and lifestyle factors in relation to
112 HCC was evaluated by means of mediation analysis. The methodological challenges
113 characterizing the analysis of large and complex metabolomic datasets are described and
114 discussed.

115 **Methods**

116 *EPIC design.* The European Prospective Investigation into Cancer and nutrition (EPIC) is a
117 large cohort established to investigate the association of diet, lifestyle and environmental
118 factors with cancer incidence and other chronic disease outcomes. Between 1992-2000, over
119 520,000 participants aged 20-85 years, were recruited from 23 centers in 10 Western
120 European countries including Denmark, France, Germany, Greece, Italy, Norway, Spain,
121 Sweden, The Netherlands and the United Kingdom [11]. The design, rationale and methods
122 of the EPIC study including information on dietary assessment methodology, blood collection
123 protocols and follow-up procedures were previously detailed [11].

124 Between 1992 and 1998, standardized lifestyle data, anthropometric measures and biological
125 samples were collected at recruitment, prior to onset of any disease [11]. Validated country-
126 specific questionnaires ensuring high compliance were used to measure diet over the previous
127 12 months [12]. Blood samples are stored at the International Agency for Research on Cancer
128 (IARC, Lyon, France) in -196°C liquid nitrogen for all countries, exceptions being Denmark
129 (nitrogen vapour, -150°C) and Sweden (freezers, -80°C).

130 *The nested case-control study.* The present study focused on data with available sera samples
131 from a nested case-control study in EPIC on hepatocellular carcinoma (HCC) [13]. Cases of
132 HCC were identified from all participating EPIC centres except for Norway and France
133 (n=117) from recruitment (1993-1998) up to 2007. Two controls (n=232) were selected for
134 each case from all cohort members alive and free of cancer (except non-melanoma skin
135 cancer) by incidence-density sampling and were matched on age at blood collection (± 1 year),
136 sex, study centre, date (± 2 months), time of the day at blood collection (± 3 hours) and fasting
137 status at blood collection (<3, 3-6, >6 hours); among women, additional matching criteria
138 included menopausal status (pre-, peri-, post-menopausal) and hormone replacement therapy
139 (HRT) use at time of blood collection (yes/no). In the present study, cases and controls were
140 both included in the analyses as the subjects were all cancer-free at blood collection. Out of
141 the total 349 subjects, 7 subjects (3 cases and 4 controls) had too little serum volume for
142 NMR spectral acquisition with sufficient sensitivity; 6 additional control subjects were
143 excluded following the exclusion of their corresponding case subject. The final analysis
144 included 114 HCC cases and 222 matched controls of which 108 case-control sets with two
145 matched control subjects and 6 sets with one matched control subject.

146 *NMR spectra acquisition.* Sera were processed using standard procedure for ^1H NMR
147 metabolic measurement and profiling protocols [14]. Details on the sera sample preparation as
148 well as NMR data acquisition and processing have been described elsewhere [15]. In brief,

149 each spectrum was reduced to 8,500 bins of 0.001 ppm width over the chemical shift range of
150 0.5 to 9 ppm. Spectra were normalized to total intensity, centred and Pareto scaled, and
151 additionally normalized for batch-effects using the batch profiling calibration method [16].
152 After removal of the structured noise (characterized by a specific mean and standard
153 deviation) located in a well-known noise region (8.5-9ppm) and variables with identical
154 characteristics, the statistical recoupling of variables (SRV) [17], a bucketing procedure, was
155 applied to the metabolomic spectra. The SRV procedure identifies clusters of variables with
156 respect to the ratio of covariance and correlation between consecutive variables along the
157 chemical shift axis, allowing the restauration of the spectral dependency and the recovery of
158 complex NMR signals corresponding to potential physical, chemical or biological entities.
159 More details on the SRV procedure are available in the **Mathematical Appendix**. This
160 permitted a reduction of the number of NMR variables from 8,500 bins to 285 clusters of
161 variables corresponding to reconstructed peak entities which constituted the Y-set of
162 metabolic variables. All steps to obtain the data were done without knowledge of the case-
163 control status of the subjects. Quality control (QC) samples were included to ensure
164 reproducibility of the NMR data acquisition.

165 *Metabolite identification.* The assignment of NMR signals observed in the ¹H one-
166 dimensional fingerprints to metabolites has been achieved by the analysis of additional 2D
167 NMR experiments ¹H-¹³C HSQC and ¹H-¹H TOCSY obtained on a subset of representative
168 samples (one control and one case). The measured chemical shifts were compared to
169 reference shifts of pure compounds using HMDB [18], MMCD [19] and ChenomX,
170 (ChenomX NMR suite, ChenomxInc, Edmonton, Canada) databases.

171 *Lifestyle variables.* The predictors (what will be referred to later on as the X-set) included 13
172 dietary variables from main EPIC food groups compiled from validated country-specific food
173 frequency questionnaires (FFQ) [11,20] (potatoes and other tubers; vegetables; legumes;

174 fruits, nuts and seeds; dairy products; cereal and cereal products; meat and meat products; fish
175 and shellfish; egg and egg products; fat; sugar and confectionary; cakes and biscuits; non-
176 alcoholic beverages), alcohol average lifetime intake (continuous, g/day), anthropometric
177 measures including body mass index (continuous, kg/m²) and height (continuous, cm) that
178 were measured by trained interviewers in the majority of participants [11], highest level of
179 education achieved (categorical: none or primary school completed, technical/professional
180 school, secondary school, longer education (incl. university degree), unspecified), smoking
181 status (categorical: never, former, current smoker, unknown), a measure of physical activity
182 (continuous, metabolic equivalents of task (MET)/h), hepatitis status (yes/no, from biomarker
183 measures of HBV and HCV seropositivity [ARCHITECT HBsAg and anti-HCV
184 chemiluminescent microparticle immunoassays; Abbott Diagnostics, France]) and baseline
185 self-reported diabetes status (yes/no). Descriptive information on these variables can be found
186 in **Supplementary table 1**.

187 *Statistical analysis*

188 *PC-PR2 analysis.* Principal component partial R-square (PC-PR2) was primarily used to
189 identify and quantify sources of systematic variability within metabolomic data [15]. PC-PR2
190 combines aspects of principal component analysis (PCA) and the R²_{partial} statistic in multiple
191 linear regression, and allows for (some) inter-correlation between the explanatory variables
192 under scrutiny [15]. In short, PCA is performed on the 285 clusters of ¹H NMR variables and
193 a number of components is retained explaining an amount of total variability above a
194 designated threshold (here, 80%). Then, multiple linear regression models are fitted where
195 each component's variability is explained in terms of relevant covariates, e.g. specific
196 characteristics of samples like country of origin, smoking status, laboratory treatment, etc. For
197 each given component, the R²_{partial} statistic is computed for all covariates, quantifying the
198 amount of variability each independent variable explains, conditional on all other covariates

199 included in the model. Finally, an overall R^2_{partial} is calculated as a weighted average for every
200 covariate, using the eigenvalues as components' weights. Mathematical details pertaining to
201 the PC-PR2 method are described elsewhere [15].

202 In this study, PC-PR2 was applied to the 285 clusters of NMR variables, whereas the
203 explanatory variables examined for systematic variability were NMR batch, country of origin,
204 sex, age at blood collection, serum clot contact time (centrifugation at the day of blood
205 collection d, or the following day, d+1), length of freezing time (≤ 15 vs. >15 years), and
206 fasting status at blood collection (< 3 , $3-6$, > 6 hours). With the similar motivation of
207 identifying sources of variability within lifestyle data, a similar PC-PR2 analysis was applied
208 to the 21 lifestyle factors; the examined covariates for systematic variability were country of
209 origin, sex and age at recruitment. For both metabolomics and lifestyle data, residuals on the
210 variable accounting for most variability, identified through PC-PR2 analyses, were computed
211 in a series of univariate linear regression models [21] and were used in the subsequent PLS.

212 *PLS analysis.* A PLS model was used to relate lifestyle variables to metabolomic profiles.
213 PLS is a multivariate technique that generalizes features of PCA and multiple linear
214 regression. PLS iteratively extracts linear combinations of, in turn, predictors (the X-set) and
215 responses (the Y-set), which in this study, were lifestyle variables and metabolomic profiles,
216 respectively. First, components or latent factors are extracted allowing a simultaneous
217 decomposition of the X- and Y-sets, in order to maximize their covariance [22]. The factors
218 extracted from the predictors' set are orthogonal. Computational details of PLS are described
219 in the **Mathematical Appendix**. As a standard step for the PLS algorithm, the X- and Y-sets
220 were centered and standardized for the analysis and a simple expectation-maximization (EM)
221 algorithm, adapted from the PLS kernel algorithm [23,24], was used to compute covariance
222 matrices when missing values were present in the lifestyle data. This was done as follows: a
223 first pass of PLS was computed filling in the missing values by the average of the non-

224 missing values for each corresponding variable. A second pass was then performed whereby
225 the missing data were assigned their predicted values based on the first model, and the PLS
226 regression is recomputed.

227 Then, a seven-fold cross validation analysis was carried out to select the number h of
228 significant PLS factors to retain [8] (see **Mathematical appendix**). This was achieved by
229 splitting the data into seven groups of observations. In turn, each group of observations was
230 considered as the test set, whilst the other six were the training sets, used to perform PLS
231 analysis. A measure of PLS performance was determined for each step through the predicted
232 residual sum of squares (PRESS) statistic, whereby the predicted values in the test set, the \hat{Y}_h
233 matrix, based on the X-components estimated through the model in the training set, were
234 compared to the observed responses, the Y matrix. This comparison is quantified by the
235 squared Euclidean distance between these two matrices. In turn for an increasing number h of
236 components, the process is iterated seven times, until each group of observations serves as a
237 test set. Eventually, the number h of selected PLS factors is the one minimizing the PRESS
238 statistic.

239 For each PLS factor, loadings were computed for the lifestyle (X-set) and the NMR (Y-set)
240 variables. The loadings, i.e. coefficients quantifying the contribution of each original variable
241 to the PLS factor, were used to characterize the various factors. As the analysis involved
242 many variables in the X-set and, particularly, in the Y-set, the interpretation focused primarily
243 on variables with loading values lower than the 10th percentile and larger than the 90th
244 percentile for the X variables, and lower than the 5th and larger than the 95th percentiles for the
245 Y variables, that were deemed the most significant contributors to the PLS factor.

246 *Logistic regression analysis.* Last, scores of each PLS factor were related to HCC risk in
247 conditional logistic regression models to compute HCC odds ratios (ORs) and associated 95%
248 confidence intervals (95% CI) where ORs express the change in HCC risk associated to one

249 standard deviation (1-SD) increase in the score. Models were adjusted for C-reactive protein
250 concentration, alpha-fetoprotein concentration and for a composite score indicative of liver
251 damage. The score summarizes the number of abnormal values of circulating enzymes
252 measured in the hepatic tissue in six liver function tests (alanine aminotransferase >55 U/L,
253 aspartate aminotransferase >34 U/L, gamma-glutamyltransferase: men>64 U/L and
254 women>36 U/L, alkaline phosphatase >150 U/L, albumin<35 g/L, total bilirubin>20.5
255 $\mu\text{mol/L}$; cut-points were provided by the clinical biochemistry laboratory that conducted the
256 analyses and were based on assay specifications) [25]. These biomarkers were measured on
257 the ARCHITECT c Systems™ and the AEROSET System (Abbott Diagnostics) using
258 standard protocols. Laboratory analyses were performed at the Centre de Biologie République
259 laboratory, Lyon, France. These adjustments were deemed necessary to address potential
260 confounding stemming from metabolic disorders, inflammation or underlying liver
261 dysfunction [25–28]. Adjustments for total dietary fibre, vitamin D, calcium and iron intakes
262 (continuous) were evaluated but not retained in the final models for lack of confounding
263 exerted by these variables. The receiver operating characteristic (ROC) curve and the
264 associated area under the curve (AUC) were determined from conditional logistic regressions
265 to evaluate the predictive performance of PLS models. AUC values were computed for
266 conditional logistic models including progressively the PLS scores, separately for lifestyle
267 and metabolomic factors (as shown in **Table 4**, column 1). The sensitivity, specificity and
268 accuracy were calculated for a cut-off point, selected as the minimal distance between the
269 ROC curve and the upper left corner of the diagram [29,30]. The corrected positive predictive
270 value (PPV), taking into account the nested case-control design [31,32] was computed by
271 including the prevalence of HCC in the EPIC population ($\pi = 0.0004$), computed over a 7-year
272 period (1992-2010) where 191 HCC cases were ascertained from a total of 477,206
273 participants included for case identification after relevant exclusions [33]. The AUC

274 unavoidably increases with the number of covariates added to the conditional logistic model.
275 To address this issue, a resampling scheme was devised to compute an objective/ unbiased
276 estimate of the AUC, inspired by the work of Uno et al [34]. For each one of the 1000 drawn
277 bootstrap samples, a 10-fold cross-validation was performed, repeated ten times to remove
278 variation due to random partitioning of data and to yield more stable estimates. The predicted
279 values from each of the conditional logistic models in the training set were used to derive
280 AUC values in the test set. The 2.5th and 97.5th percentile values made up the 95% confidence
281 intervals.

282 *Sensitivity analyses.* A sensitivity analysis was performed by running PLS on data excluding
283 sets where cases were diagnosed within the first two years of follow-up. The model was
284 conducted on 271 observations (92 cases, 179 controls), to investigate the performance of the
285 PLS model, ruling out potential reverse causation. The metabolomic profiles of HCC cases
286 diagnosed within two years from enrolment could reflect the presence of the tumour rather
287 than informing about tumour aetiology. The variable importance in the projection (VIP)
288 statistic was used to facilitate the comparison of the sensitivity analysis with the main
289 analysis. The VIP expresses the explanatory power of a predictor variable X across all
290 response variables Y (see **Mathematical Appendix**).

291 *Mediation analysis.* The mediating role of the Y-scores in the association between lifestyle
292 profiles and HCC risk was assessed. Separately for each extracted combination of lifestyle
293 and metabolomic PLS factors, mediation analyses were performed with the ‘paramed’ Stata
294 function that allows for exposure-mediator interaction based on Valeri and VanderWeele’s
295 work [35]. Briefly, mediation was computed using a Baron and Kenny approach adapted to
296 dichotomous outcomes [36], where two models were specified. In the mediator model, the
297 mediator (the Y-score) was linearly regressed on the exposure (the X-score), while in the
298 outcome model the exposure (X-score) and the mediator (Y-score) were related to the HCC

299 indicator in unconditional logistic regressions. Both models accounted for the concentration
300 of C-reactive protein, alpha-fetoprotein and the composite score of liver damage, and
301 additionally accommodated the other extracted metabolic profiles (Y-scores) to control for
302 mediator-outcome confounders that may occur when estimating the Natural Indirect Effect
303 (NIE) [35]. As the outcome (HCC) is rare, direct and indirect effects can be estimated taking
304 into account the case-control design. This is done by using the same formulas for the effects,
305 while running the mediator regression only for the controls [36]. As mediation packages do
306 not yet accommodate conditional logistic models, the outcome and the mediator models,
307 which were accommodated in unconditional logistic regressions, were adjusted for center and
308 age at blood collection for sake of consistency with previous steps of the analysis.

309 Statistical analyses were performed using R [37] and SAS [38] in general, with the following
310 packages for specific purposes: PROC PLS in SAS 9.4 for PLS analyses, ‘paramed’ in Stata
311 12 [39] for mediation analyses, ‘OptimalCutpoints’ in R for ROC-related assessments.

312 The different steps of the analytical framework developed in this study to model the MITM
313 are presented in **Figure 1**.

314 **Results**

315 In the PC-PR2 analyses, a total of 17 and 14 principal components were retained to
316 explain an amount of total variability exceeding 80% in metabolomics and lifestyle data
317 respectively. **Figure 2** shows that the ensemble of explanatory variables accounted for 19.4%
318 and 26.7% of total variance, respectively in metabolomics and lifestyle data, of which the
319 highest contributor was ‘country of origin’ with consistently 8% and 22%. -Major sources of
320 variation in the Y-set displaying large R^2_{partial} value were country of origin (8.3%), NMR
321 batch (4.0%) and fasting status at blood collection (1.6%). In the X-set, country of origin
322 (22.6%) and sex (5.1%) showed the highest contributions. As PC-PR2 analysis showed that

323 ~~'country of origin' accounted for about 8% of the variability within the metabolomic data, and~~
324 ~~22% in the lifestyle variables, the~~ PLS analysis was carried controlling for this variable.

325 After a seven-fold cross-validation, three PLS factors were retained accounting for
326 21.7% and 8.5% of the overall variability observed in predictor and response variables,
327 respectively (**Table 1**). Lifestyle variables and clusters of NMR variables contributing highly
328 to PLS factors were identified using factor loading values (**Table 2**). The first PLS factor was
329 predominantly positively associated with dairy products and cakes and biscuits intake, while
330 lifetime alcohol intake, smoking status and diabetes displayed negative loadings for this
331 lifestyle component (**Table 2**). On the same PLS factor, signals mainly associated with
332 glucose and bonds of lipids with negative loading values, and with aspartate, glutamine and
333 lysine with positive loadings emerged on the metabolomic profile (**Table 2**). Lifestyle
334 variables characterizing the second PLS factor included cereal products, height and education
335 level with negative loadings, and hepatitis with positive loadings. The metabolic signature
336 included NMR variables with positive loadings associated with aromatic amino acids
337 (phenylalanine, tyrosine) and glucose; and those with negative loadings associated mainly
338 with bonds of lipids, threonine and mannose (**Table 2**). The third PLS factor had a lifestyle
339 pattern outlining intake of vegetables (high negative loadings values), lifetime alcohol
340 consumption, smoking, and hepatitis infection (positive loadings). Its counterpart NMR
341 pattern highlighted signals of glucose and aspartate, with high negative loadings, along with
342 signals of ethanol, myo-inositol, proline and glutamate as prominent metabolites with positive
343 loadings (**Table 2**).

344 Conditional logistic regression models relating HCC risk with the X- and Y-scores are
345 shown in **Table 3**. The first PLS factor was associated to a non-significant decreased HCC
346 risk (23% and 4% in the X- and Y-scores respectively), while the second and third factors
347 were associated to a statistically significant increased HCC risk (54% and 11%; and 37% and

348 22% respectively). Results for the ROC curves parameters are reported in **Table 4**, including
349 AUC, sensitivity, specificity, accuracy and PPV for different combinations of the X- and Y-
350 scores. The AUC of the X-scores and Y-scores for all 3 PLS factors, adjusted for C-reactive
351 protein concentration, alpha-fetoprotein concentration and the score of liver damage, was
352 respectively 0.859 and 0.853. An increase in the resampled cross-validated AUC values was
353 also observed for all three X- and Y-scores, albeit smaller, with respectively 0.836 and 0.827.
354 Results from the sensitivity analysis conducted on data excluding sets where cases were
355 diagnosed within the first two years of follow-up, showed similarities in terms of lifestyle
356 variables' and metabolites' loadings on the PLS factors (**Supplementary Table 2**). Notable
357 differences pertained to the identification of new signals for the first PLS factor including
358 ethanol, histidine and an unknown compound. On the second lifestyle factor, BMI (positive
359 loadings) replaced education level (negative loadings) while the reflected metabolomic profile
360 was comparable to its counterpart from the main analysis (**Supplementary Table 2**). On the
361 third factor, smoking status and hepatitis (positive loadings) were replaced by sugar and
362 confectionary intake (negative loadings); signals contributing to the associated metabolic
363 profile remained the same but the direction of the association was inversed as loadings had
364 opposite signs as compared to the counterpart PLS factor of the main model (**Supplementary**
365 **Table 2**). Corresponding ORs from conditional logistic regression models relating the X- and
366 Y-scores to HCC risk are available in **Table 5**. The scores showed a statistically significant
367 association in the second factor for both sets and in the third factor for the Y-set. ROC-
368 associated statistics for different models are presented in **Supplementary Table 3**. The VIP
369 plot (**Figure 3**) displayed the results for the importance of the lifestyle variables in the
370 prediction of the Y-set computed for the main PLS model performed including all subjects
371 (panel **A**) and for the sensitivity model (panel **B**). The results suggested a potential gain in
372 stability as prominent lifestyle variables for prediction were maintained

373 (hepatitis/diabetes/cakes and biscuits), the magnitude of the VIP was improved for some
374 (fat/lifetime alcohol intake) and less emphasis was put on others (BMI/physical activity).
375 Finally, the natural indirect effect was assessed in the mediation analyses and the results are
376 presented in **Table 6**. Overall, there was limited evidence that metabolomic signals mediated
377 the association between lifestyle components and HCC risk in the first PLS factor. Evidence
378 of a significant mediated effect by the Y-scores was found in the second and third PLS factors
379 when models were adjusted for exposure-mediator interaction (**Table 6**).

380 **Discussion**

381 In this work, an analytical strategy based on PLS analysis was conceived to extract
382 relevant information from sets of lifestyle and NMR metabolomic variables, and to relate the
383 resulting components to the risk of disease. This offered a way to implement the MITM
384 approach [6] in a nested case-control study on HCC within the EPIC study. MITM has been
385 suggested as a way to link specific putative metabolites to lifestyle exposures and disease
386 outcomes, thus leading to the identification of potential intermediate biomarkers [6].

387 An implementation of MITM was previously carried out in a nested case-control study
388 in the Turin sub-cohort of EPIC [7] based on prospectively collected plasma samples from a
389 pilot study on colon and breast cancers. In their work, a list of intermediate markers was
390 identified by an in-parallel evaluation of the relationships between untargeted ¹H NMR
391 profiles with dietary exposures and risk of colon and breast cancers using correlation analysis
392 and logistic regression. In our study, a different analytical framework was developed, largely
393 exploiting features of PLS analysis, a multivariate technique that iteratively extracts
394 components capturing co-variability in sets of predictors and response variables [8,40]. A set
395 of lifestyle predictor variables were related to NMR responses. In a second step, PLS
396 predictors' and responses' scores were linked to the risk of HCC.

397 Another sensitive issue in this analysis was the choice of lifestyle variables. Two
398 disease-indicator variables reflecting environmental exposures, diabetes and hepatitis, were
399 included in the set of predictors, as they turned out to have an important role in the
400 characterization of metabolomic signatures. In addition, diabetes is the main metabolic risk
401 factor for HCC alongside with fatty liver disease [41,42], and chronic infection with hepatitis
402 B (HBV) and particularly hepatitis C (HCV) viruses were classified as class I carcinogens for
403 HCC by IARC [43].

404 Other relevant biomarkers were not part of the list of predictors in PLS analysis, but were
405 controlled for in logistic regression models. This included C-reactive protein, alpha-
406 fetoprotein, and a score for liver damage, an index of different circulating enzymes measured
407 in the hepatic tissue indicating potential underlying liver function impairment [25]. The alpha-
408 fetoprotein was ~~not~~ included as an adjustment factor in the analyses not because of its
409 established part as a serum marker for HCC diagnosis [26,44], but rather to account for it as a
410 potential confounder that may cloud the relation between scores and HCC, both in conditional
411 logistic regressions and in mediation analyses.

412 Similarly to other multivariate techniques, a key aspect of PLS analysis is the choice
413 of the number of factors to retain, in an effort of exhaustively summarizing data variability
414 through a limited number of factors. Based on a seven-fold cross-validation, three linear
415 combinations of variables were extracted in this work. A challenging aspect of this analysis is
416 the interpretation of these factors, with respect to lifestyle and metabolomic variables. A
417 subjective criterion based on the distribution of loading values was used throughout. The
418 variables displaying the most extreme loading values (in absolute terms) were the ones
419 characterizing each factor.

420 The first lifestyle factor highlighted a healthy pattern with negative loadings for
421 diabetes status, smoking status and lifetime alcohol intake, and was not associated to HCC

422 risk, similarly to its metabolomics counterpart. The lifestyle component of the second PLS
423 factor, was reflective of a lifestyle pattern reflective of “higher-risk exposures”, and was
424 related to a significant 54% increase in HCC risk. Likewise, its associated metabolic
425 component displayed a significant HCC risk augmentation by 11%. The lifestyle component
426 of the third PLS factor described participants with lower vegetables intake, elevated lifetime
427 alcohol consumption, more likely to be ever smokers and hepatitis positive; one standard
428 deviation increase of this component was associated to a statistically significant 37% increase
429 in HCC risk. Similarly, a 22% significant increase in HCC risk was observed for its metabolic
430 counterpart, characterized by positive signals of ethanol and myo-inositol, and displayed
431 negative loadings for glucose.

432 The MITM is captured by the rationale of PLS analysis, in the sense that each set of lifestyle
433 profiles and metabolic signatures of the extracted PLS factors mirrored one another. In
434 addition, mediation was observed for the second and third PLS factors, whereby the
435 metabolomic component mediated the relation between the lifestyle component and HCC, for
436 which statistically significant associations with HCC risk were estimated, emphasising the
437 presence of a MITM. Mediation analysis relies on the assumption that there is no mediator-
438 outcome confounder that is affected by the exposure [35]. In our study C-reactive protein,
439 alpha-fetoprotein and liver damage score were weakly correlated to lifestyle factor score, thus
440 introducing potential bias in the estimation of direct and indirect effects in our mediation
441 analysis. Additionally, a number of background confounders (mediator-outcome and
442 exposure-outcome confounders) were present that we have tried to control for, either by
443 adjustments or by accounting for potential interactions, however some degree of bias can
444 remain and caution should be employed when interpreting the results.

445 The predictive performance of PLS factors in relation to HCC occurrence was evaluated
446 through an analysis of AUC values. The performance of the model ~~was~~ improved

447 progressively, ~~with all 3 X- and Y-scores added; and reached an AUC of 0.859 with all 3 X-~~
448 ~~scores and 0.853 with all 3 Y-scores, with adjustment for concentrations of C-reactive and~~
449 ~~alpha-fetoprotein, and for the liver damage score. A~~fter a bootstrapped cross-validation, the
450 AUC estimates were lower ~~with respectively 0.836 and 0.827, but~~ the increase in the
451 performance was nevertheless present. The ROC methodology allows estimation of PPV,
452 which expresses the risk of disease after a positive test [45]. In a setting with low HCC
453 prevalence ($\pi=0.0004$), in line with Western populations [46], extremely low PPV estimates
454 were observed. In the absence of a very specific test, many positive tests arise from disease-
455 free individuals [45], thus leading to a dilution of PPV.

456 A sensitivity analysis was carried out excluding the first two years of follow-up, but results
457 were virtually unchanged, both in terms of relative risk estimates in logistic regression
458 models, and of percentage of variability explained in PLS analysis. These findings suggest
459 that reverse causation bias, if present, was minimal.

460 This study had the ambition of integrating in the same analytical framework study
461 participants' lifestyle characteristics with a large number of NMR metabolic profiles. These
462 data pose a number of methodological challenges due to their size and the complexity of
463 exhaustively capturing and interpreting the biological processes they reflect. To address these
464 issues, techniques involving multivariate statistics have been progressively revived in the
465 recent years [2]. Epidemiologic evaluations of metabolomic data frequently combined PLS
466 with discriminant analysis, such as PLS-DA or O-PLS-DA. The main objective of these
467 methods is to identify a series of metabolomic features distinguishing between two very
468 distinct groups of study participants [47,48]. In such strategies, only one set of variables is
469 multi-dimensional and the response is one variable only. Similar multivariate techniques for
470 pattern extraction, belonging to the family of regression methods, include reduced rank
471 regression. This multivariate method relates an ensemble of response variables to a set of

472 predictor variables where the estimated matrix of the regression coefficients is of reduced
473 rank [49–51]. In addition, canonical correlation analysis (CCA) [52] is a method applied to
474 identify the optimum structure or dimensionality of each variable set that maximizes the
475 relationship between two sets of multi-dimensional variables. The main difference between
476 CCA and PLS regression is that CCA maximizes the correlation between the two new
477 dimensions, i.e. extracted factors, whereas PLS maximizes their covariance. PLS can be
478 considered as a trade-off between CCA and PCA, since maximizing the covariance
479 corresponds to maximizing the product of the correlation and standard deviation, given that
480 $\text{cov}(X,Y)=\text{cor}(X,Y)*\text{SD}(X)*\text{SD}(Y)$.

481 Untargeted NMR was used in this work to acquire metabolomic signals. Prior to PLS
482 analysis, a bucketing procedure, the statistical recoupling of variables (SRV) [17,53], was
483 applied to reduce the number of NMR variables to 285 clusters. This was done by aggregating
484 consecutive NMR bins based on their covariance to correlation ratio, ~~thus reconstructing~~
485 ~~peak entities. Neighbouring variables were then merged into clusters, to recover NMR~~
486 ~~multiplets, corresponding to NMR variables of interest.~~ This allows the identification of
487 informative components of the spectra, thus acting as an efficient noise-removing filter.
488 Subsequently the annotation effort remains challenging, for a number of reasons. The
489 majority of published metabolomics studies often identified a limited number of metabolites
490 at a time [54], and the Human Metabolome Database (HMDB) and other related resources
491 [18,55], that offer richly annotated information continuously increasing the metabolite
492 coverage for users, are mostly exploited through time consuming interactive procedures. In
493 addition, individual metabolites often overlap in NMR signals, which can hinder
494 ~~interpretation of the annotated metabolic profiles annotations. Untargeted NMR approaches~~
495 ~~are useful for the identification of metabolites of moderate abundance; however they may~~
496 ~~miss low abundance metabolites due to the intrinsically low NMR sensitivity.~~ These

497 challenges, as well as large variability in metabolite concentrations, and disentangling
498 informative signals from noise, are not specific to NMR and pertain to any type of untargeted
499 technique. Such investigations may profit from complementary targeted metabolomic
500 analytical strategies [55].

501 Throughout the different steps of this work, the scaling problem was first tackled by
502 normalizing spectra to total intensity. NMR data were also centered and Pareto-scaled,
503 together with correction for potential batch effects [16]. The PC-PR2 method offered a way to
504 investigate major sources of systematic variability in NMR and lifestyle data [15]. The
505 variable “country of origin” emerged as the variable accounting for the largest proportion of
506 total variability, and the residual method was used to control for this variable in the following
507 steps of the analysis. While this may lead to removing regional gradients of dietary
508 variability, this step is instrumental to avoid unwanted systematic regional-specific bias in the
509 data in country-specific questionnaire assessments. In addition, technical aspects like storage
510 and handling of biological samples, fasting status at blood collection are specific to each
511 country [15]. In any case, variability due to “country of origin” is not exploited in conditional
512 logistic models, as cases and controls were also matched on center.

513 One of the limitations of this study is the restricted sample size which raises concerns
514 with regards to power to detect associations. While a larger sample size would possibly result
515 in more statistically significant findings, we used the data that was available with NMR
516 profiles measured. In this work we have developed a framework to analyse complex data
517 integrating lifestyle and metabolomics in relation to risk of disease. The approach described in
518 this study has merits but also pitfalls among which it is worth mentioning that statistical
519 methods are used repeatedly on the same set of data, notably the PLS model, the conditional
520 logistic regression, the AUC estimation and mediation analysis. To partially address this, a
521 cross-validation approach was devised for AUC estimation which involved conditional

522 logistic regression, whereby PLS was done without knowledge of the case/control status.
523 However, conditional logistic regression models and mediation analyses were implemented
524 on the same data, and our analysis did not account for this limitation. This may have led to
525 spuriously increase the nominal level of statistical significance of statistical tests.

526 **Conclusion**

527 The MITM emerged as a method for the identification of relevant biomarkers, with
528 great potential to unravel utmost important steps in the aetiology of disease. The analytical
529 strategy for MITM was developed to use all potentially informative aspects of high-
530 throughput data by integrating metabolomic, dietary and lifestyle exposures together with
531 disease indicators. While the framework was applied towards the investigation of HCC
532 determinants, it can be easily extended to similar aetiological contexts and applied to other –
533 omics settings.

534 **Funding**

535 This work was supported by the French National Cancer Institute (L'Institut National du
536 Cancer; INCA) [grant number 2009-139; PI: M. Jenab]. The coordination of the European
537 Prospective Investigation into Cancer and nutrition is financially supported by the European
538 Commission (Directorate General for Health and Consumer Affairs) and the International
539 Agency for Research on Cancer. The national cohorts are supported by Health Research Fund
540 (FIS) of the Spanish Ministry of Health RTICC 'Red Temática de Investigación Cooperativa
541 en Cáncer [Grant numbers: Rd06/0020/0091, Rd12/0036/0018], Regional Governments of
542 Andalucía, Asturias, Basque Country, Murcia [project 6236] and Navarra, Instituto de Salud
543 Carlos III, Redes de Investigación Cooperativa (RD06/0020) (Spain); the Danish Cancer
544 Society (Denmark); the Ligue Contre le Cancer, the Institut Gustave Roussy, the Mutuelle
545 Générale de l'Éducation Nationale, and the Institut National de la Santé et de la Recherche
546 Médicale (France); the Deutsche Krebshilfe, the Deutsches Krebsforschungszentrum, and the
547 Federal Ministry of Education and Research (Germany); the Hellenic Health Foundation, the
548 Stavros Niarchos Foundation, and the Hellenic Ministry of Health and Social Solidarity
549 (Greece); the Italian Association for Research on Cancer AIRC and the National Research
550 Council (Italy); the Dutch Ministry of Public Health, Welfare and Sports, the Netherlands

551 Cancer Registry, LK Research Funds, Dutch Prevention Funds, the Dutch Zorg Onderzoek
552 Nederland, the World Cancer Research Fund, and Statistics Netherlands (Netherlands);
553 European Research Council-2009-AdG 232997 and the Nordforsk, Nordic Centre of
554 Excellence program on Food, Nutrition and Health (Norway); the Swedish Cancer Society,
555 the Swedish Research Council, and the Regional Governments of Skåne and Västerbotten
556 (Sweden); Cancer Research UK, the Medical Research Council, the Stroke Association, the
557 British Heart Foundation, the Department of Health, the Food Standards Agency, and the
558 Wellcome Trust (United Kingdom). The work undertaken by N Assi was supported by by the
559 Université de Lyon I through a doctoral fellowship awarded by the EDISS doctoral school.

560 **Acknowledgments**

561 We would like to acknowledge the assistance of Dr Elodie Jobard from the ISA-CRMN in
562 obtaining the annotation of the NMR data.

563 **Reference List**

- 564 1. Nicholson, J.K., Holmes, E., and Elliott, P. (2008) The metabolome-wide association
565 study: a new look at human disease risk factors. *J. Proteome Res.*, **7**,3637–3638.
- 566 2. Chadeau-Hyam, M., Campanella, G., Jombart, T., Bottolo, L., Portengen, L., Vineis,
567 P., Liquet, B., and Vermeulen, R.C.H. (2013) Deciphering the Complex:
568 Methodological Overview of Statistical Models to Derive OMICS-Based Biomarkers.
569 *Environ. Mol. Mutagen.*, **54**,542–557.
- 570 3. Floegel, A., Wientzek, A., Bachlechner, U., Jacobs, S., Drogan, D., Prehn, C.,
571 Adamski, J., Krumsiek, J., Schulze, M.B., Pischon, T., and Boeing, H. (2014) Linking
572 diet, physical activity, cardiorespiratory fitness and obesity to serum metabolite
573 networks: findings from a population-based study. *Int. J. Obes. (Lond)*,1–9.
- 574 4. Trushina, E. and Mielke, M.M. (2013) Recent advances in the application of
575 metabolomics to Alzheimer's Disease. *Biochim. Biophys. Acta*, **1842**,1232–1239.
- 576 5. Jin, X., Yun, S.J., Jeong, P., Kim, I.Y., Kim, W.-J., and Park, S. (2014) Diagnosis of
577 bladder cancer and prediction of survival by urinary metabolomics. *Oncotarget*,
578 **5**,1635–1645.
- 579 6. Vineis, P. and Perera, F. (2007) Molecular epidemiology and biomarkers in etiologic
580 cancer research: the new in light of the old. *Cancer Epidemiol. Biomarkers Prev.*,
581 **16**,1954–1965.
- 582 7. Chadeau-Hyam, M., Athersuch, T.J., Keun, H.C., De Iorio, M., Ebbels, T.M.D., Jenab,
583 M., Sacerdote, C., Bruce, S.J., Holmes, E., and Vineis, P. (2011) Meeting-in-the-
584 middle using metabolic profiling - a strategy for the identification of intermediate
585 biomarkers in cohort studies. *Biomarkers*, **16**,83–88.
- 586 8. Tenenhaus, M. (1998) La régression PLS. Technip. Paris.
- 587 9. Mitra, V. and Metcalf, J. (2009) Metabolic functions of the liver. *Anaesth. Intensive*
588 *Care Med.*, **10**,334–335.
- 589 10. Fages, A., Duarte-Salles, T., Stepien, M., Ferrari, P., Fedirko, V., Pontoizeau, C.,
590 Trichopoulou, A., Aleksandrova, K., Tjønneland, A., Olsen, A., Clavel-Chapelon, F.,
591 Boutron-Ruault, M.-C., Severi, G., Kaaks, R., Kuhn, T., Floegel, A., Boeing, H.,
592 Lagiou, P., Bamia, C., Trichopoulos, D., Palli, D., Pala, V., Panico, S., Tumino, R.,
593 Vineis, P., Bueno-de-Mesquita, H.B., Peeters, P.H.M., Weiderpass, E., Agudo, A.,
594 Molina-Montes, E., Huerta, J.M., Ardanaz, E., Dorronsoro, M., Sjöberg, K., Ohlsson,
595 B., Khaw, K.-T., Wareham, N.J., Travis, R.C., Schmidt, J.A., Cross, A.J., Gunter, M.J.,
596 Riboli, E., Scalbert, A., Romieu, I., Elena-Herrmann, B., and Jenab, M. (2015)
597 Metabolomic Profiles of Hepatocellular Carcinoma in a European Prospective Cohort.
598 *Submitt. to Am. J. Gastroenterol.*,.
- 599 11. Riboli, E., Hunt, K.J., Slimani, N., Ferrari, P., Norat, T., Fahey, M., Charrondière,
600 U.R., Hémon, B., Casagrande, C., Vignat, J., Overvad, K., Tjønneland, A., Clavel-
601 Chapelon, F., Thiébaud, A., Wahrendorf, J., Boeing, H., Trichopoulos, D.,

Commented [NA1]: The references, format has been modified according to the Editor's request to include all authors.

- 602 Trichopoulou, A., Vineis, P., Palli, D., Bueno-De-Mesquita, H.B., Peeters, P.H.M.,
603 Lund, E., Engeset, D., González, C. a, Barricarte, A., Berglund, G., Hallmans, G., Day,
604 N.E., Key, T.J., Kaaks, R., and Saracci, R. (2002) European Prospective Investigation
605 into Cancer and Nutrition (EPIC): study populations and data collection. *Public Health*
606 *Nutr.*, **5**,1113–1124.
- 607 12. Kaaks, R., Slimani, N., and Riboli, E. (1997) Pilot Phase Studies on the Accuracy of
608 Dietary Intake Measurements in the EPIC Project : Overall Evaluation of Results.
609 **26**,26–36.
- 610 13. Trichopoulos, D., Bamia, C., Lagiou, P., Fedirko, V., Trepo, E., Jenab, M., Pischon, T.,
611 Nöthlings, U., Overved, K., Tjønneland, A., Outzen, M., Clavel-Chapelon, F., Kaaks,
612 R., Lukanova, A., Boeing, H., Aleksandrova, K., Benetou, V., Zylis, D., Palli, D., Pala,
613 V., Panico, S., Tumino, R., Sacerdote, C., Bueno-De-Mesquita, H.B., Van Kranen,
614 H.J., Peeters, P.H.M., Lund, E., Quirós, J.R., González, C. a, Sanchez Perez, M.-J.,
615 Navarro, C., Dorronsoro, M., Barricarte, A., Lindkvist, B., Regnér, S., Werner, M.,
616 Hallmans, G., Khaw, K.-T., Wareham, N., Key, T., Romieu, I., Chuang, S.-C.,
617 Murphy, N., Boffetta, P., Trichopoulou, A., and Riboli, E. (2011) Hepatocellular
618 carcinoma risk factors and disease burden in a European cohort: a nested case-control
619 study. *J. Natl. Cancer Inst.*, **103**,1686–1695.
- 620 14. Beckonert, O., Keun, H.C., Ebbels, T.M.D., Bundy, J., Holmes, E., Lindon, J.C., and
621 Nicholson, J.K. (2007) Metabolic profiling, metabolomic and metabonomic procedures
622 for NMR spectroscopy of urine, plasma, serum and tissue extracts. *Nat. Protoc.*,
623 **2**,2692–2703.
- 624 15. Fages, A., Ferrari, P., Monni, S., Dossus, L., Floegel, A., Mode, N., and Al., E. (2014)
625 Investigating sources of variability in metabolomic data in the EPIC study: the
626 Principal Component Partial R-square (PC-PR2) method. *Metabolomics*, **10**,1074–
627 1083.
- 628 16. Fages, A., Pontoizeau, C., Jobard, E., Lévy, P., Bartosch, B., and Elena-Herrmann, B.
629 (2013) Batch profiling calibration for robust NMR metabonomic data analysis. *Anal.*
630 *Bioanal. Chem.*, **405**,8819–8827.
- 631 17. Blaise, B.J., Shintu, L., Elena, B., Emsley, L., Dumas, M.-E., and Toulhoat, P. (2009)
632 Statistical Recoupling Prior to Significance Testing in Nuclear Resonance Based
633 Metabonomics. *Anal. Chem.*, **81**,6242–6251.
- 634 18. Wishart, D.S., Knox, C., Guo, A.C., Eisner, R., Young, N., Gautam, B., Hau, D.D.,
635 Psychogios, N., Dong, E., Bouatra, S., Mandal, R., Sinelnikov, I., Xia, J., Jia, L., Cruz,
636 J. a, Lim, E., Sobsey, C. a, Shrivastava, S., Huang, P., Liu, P., Fang, L., Peng, J.,
637 Fradette, R., Cheng, D., Tzur, D., Clements, M., Lewis, A., De Souza, A., Zuniga, A.,
638 Dawe, M., Xiong, Y., Clive, D., Greiner, R., Nazyrova, A., Shaykhtudinov, R., Li, L.,
639 Vogel, H.J., and Forsythe, I. (2009) HMDB: a knowledgebase for the human
640 metabolome. *Nucleic Acids Res.*, **37**,D603–D610.
- 641 19. Cui, Q., Lewis, I.A., Hegeman, A.D., Anderson, M.E., Li, J., Schulte, C.F., Westler,
642 W.M., Eghbalnia, H.R., Sussman, M.R., and Markley, J.L. (2008) Metabolite

- 643 identification via the Madison Metabolomics Consortium Database. *Nat. Biotechnol.*,
644 **26**,162–164.
- 645 20. Slimani, N., Deharveng, G., Unwin, I., Southgate, D.A., Viignat, J., Skeie, G., Salvini,
646 S., Parpinel, M., Moller, A., and et al. (2007) The EPIC nutrient database project
647 (ENDB): a first attempt to standardize nutrient databases across the 10 European
648 countries participating in the EPIC study - DTU Orbit. *Eur. J. Clin. Nutr.*, **61**,1037–
649 1056.
- 650 21. Kleinbaum, D.G., Kupper, L.K., and Muller, K.E. (1987) Applied regression analysis
651 and other multivariable methods. Belmont, CA: Duxbury Press.
- 652 22. Wold, S., Sjoström, M., and Ericksson, L. (2001) PLS-regression : a basic tool of
653 chemometrics. *Chemom. Intell. Lab. Syst.*, **58**,109–130.
- 654 23. Rannar, S., Geladi, P., Lindgren, F., and Wold, S. (1995) A PLS kernel algorithm for
655 data sets with many variables and few objects. Part II: cross-validation, missing data
656 and examples. *J. Chemom.*, **9**,459–470.
- 657 24. Bastien, P. (2008) Régression PLS et Données Censurées. Conservatoire National des
658 Arts et Métiers - CNAM.
- 659 25. Fedirko, V., Trichopolou, A., Bamia, C., Duarte-Salles, T., Trepo, E., Aleksandrova,
660 K., Nöthlings, U., Lukanova, A., Lagiou, P., Boffetta, P., Trichopoulos, D., Katzke, V.
661 a, Overvad, K., Tjønneland, A., Hansen, L., Boutron-Ruault, M.C., Fagherazzi, G.,
662 Bastide, N., Panico, S., Grioni, S., Vineis, P., Palli, D., Tumino, R., Bueno-de-
663 Mesquita, H.B., Peeters, P.H., Skeie, G., Engeset, D., Parr, C.L., Jakszyn, P., Sánchez,
664 M.J., Barricarte, A., Amiano, P., Chirlaque, M., Quirós, J.R., Sund, M., Werner, M.,
665 Sonestedt, E., Ericson, U., Key, T.J., Khaw, K.T., Ferrari, P., Romieu, I., Riboli, E.,
666 and Jenab, M. (2013) Consumption of fish and meats and risk of hepatocellular
667 carcinoma: the European Prospective Investigation into Cancer and Nutrition (EPIC).
668 *Ann. Oncol.*, **24**,2166–2173.
- 669 26. Akuta, N., Suzuki, F., Kobayashi, M., Hara, T., Sezaki, H., Suzuki, Y., Hosaka, T.,
670 Kobayashi, M., Saitoh, S., Ikeda, K., and Kumada, H. (2014) Correlation Between
671 Hepatitis B Virus Surface antigen Level and Alpha-Fetoprotein in Patients Free of
672 Hepatocellular Carcinoma or Severe Hepatitis. *J. Med. Virol.*, **86**,131–138.
- 673 27. Kanazir, M., Boricic, I., Delic, D., Tepavcevic, D.K., Knezevic, A., Jovanovic, T., and
674 Pekmezovic, T. (2010) Risk factors for hepatocellular carcinoma: A case-control study
675 in Belgrade (Serbia). *Tumori*, **96**,911–917.
- 676 28. Zheng, Z., Zhou, L., Gao, S., Yang, Z., Yao, J., and Zheng, S. (2013) Prognostic role of
677 C-reactive protein in hepatocellular carcinoma: A systematic review and meta-analysis.
678 *Int. J. Med. Sci.*, **10**,653–664.
- 679 29. Metz, C.D. (1978) Basic Principles of ROC Analysis. *Semin. Nucl. Med.*, **8**,283–298.

- 680 30. Vermont, J., Bosson, J.L., François, P., Robert, C., Rueff, A., and Demongeot, J.
681 (1991) Strategies for graphical threshold determination. *Comput. Methods Programs*
682 *Biomed.*, **35**,141–150.
- 683 31. Biesheuvel, C.J., Vergouwe, Y., Oudega, R., Hoes, A.W., Grobbee, D.E., and Moons,
684 K.G.M. (2008) Advantages of the nested case-control design in diagnostic research.
685 *BMC Med. Res. Methodol.*, **8**,48.
- 686 32. Van Zaane, B., Vergouwe, Y., Donders, a R.T., and Moons, K.G.M. (2012)
687 Comparison of approaches to estimate confidence intervals of post-test probabilities of
688 diagnostic test results in a nested case-control study. *BMC Med. Res. Methodol.*,
689 **12**,166.
- 690 33. Stepien, M., Duarte-Salles, T., Fedirko, V., Floegel, A., Kumar-Barupal, D., Rinaldi,
691 S., Achaintre, D., Assi, N., Tjønneland, A., Overvad, K., Bastide, N., Boutron-Ruault,
692 M.-C., Severi, G., Kuhn, T., Kaaks, R., Aleksandrova, K., Boeing, H., Trichopoulou,
693 A., Bamia, C., Lagiou, P., Saieva, C., Agnoli, C., Panico, S., Tumino, R., Naccarati, A.,
694 Bueno-de-Mesquita, H.B., Peeters, P.H., Weiderpass, E., Quirós, J.R., Agudo, A.,
695 Sanchez, M.-J., Dorronsoro, M., Gavrilu, D., Barricarte, A., Ohlsson, B., Sjöberg, K.,
696 Werner, M., Sund, M., Wareham, N., Khaw, K.-T., Travis, R.C., Schmidt, J.A., Gunter,
697 M., Cross, A.J., Vineis, P., Romieu, I., Scalbert, A., and Jenab, M. (2015) Alteration of
698 Amino Acid and Biogenic Amine Metabolism in Hepatobiliary Cancers: Findings from
699 a Prospective Cohort Study. *Submitt. to Int. J. Cancer*,.
- 700 34. Uno, H., Cai, T., Pencina, M.J., D'Agostino, R.B., and Wei, L.J. (2011) On the C-
701 statistics for Evaluating Overall Adequacy of Risk Prediction Procedures with
702 Censored Survival Data. *Stat. Med.*, **30**,1105–1117.
- 703 35. Valeri, L. and Vanderweele, T.J. (2013) Mediation analysis allowing for exposure-
704 mediator interactions and causal interpretation: theoretical assumptions and
705 implementation with SAS and SPSS macros. *Psychol. Methods*, **18**,137–150.
- 706 36. Vanderweele, T.J. and Vansteelandt, S. (2010) Odds ratios for mediation analysis for a
707 dichotomous outcome. *Am. J. Epidemiol.*, **172**,1339–1348.
- 708 37. R Foundation for Statistical Computing and R Core Team. (2013) R: A language and
709 environment for statistical computing.
- 710 38. SAS Institute Inc. and Cary, N. (2012) Base SAS® 9.4 Procedures Guide.
- 711 39. College Station TX : StataCorp. (2011) Stata Statistical Software: Release 12.
- 712 40. Abdi, H. (2010) Partial least squares regression and projection on latent structure
713 regression (PLS Regression). *Wiley Interdiscip. Rev. Comput. Stat.*, **2**,97–106.
- 714 41. Yang, W.-S., Va, P., Bray, F., Gao, S., Gao, J., Li, H.-L., and Xiang, Y.-B. (2011) The
715 role of pre-existing diabetes mellitus on hepatocellular carcinoma occurrence and
716 prognosis: a meta-analysis of prospective cohort studies. *PLoS One*, **6**,e27326.

- 717 42. Gomaa, A.-I. (2008) Hepatocellular carcinoma: Epidemiology, risk factors and
718 pathogenesis. *World J. Gastroenterol.*, **14**,4300–4308.
- 719 43. Cogliano, V.J., Baan, R., Straif, K., Grosse, Y., Lauby-Secretan, B., El Ghissassi, F.,
720 Bouvard, V., Benbrahim-Tallaa, L., Guha, N., Freeman, C., Galichet, L., and Wild,
721 C.P. (2011) Preventable exposures associated with human cancers. *J. Natl. Cancer*
722 *Inst.*, **103**,1827–1839.
- 723 44. Bialecki, E.S. and Di Bisceglie, A.M. (2005) Diagnosis of hepatocellular carcinoma.
724 *HPB (Oxford)*, **7**,26–34.
- 725 45. Wentzensen, N. and Wacholder, S. (2013) From differences in means between cases
726 and controls to risk stratification: a business plan for biomarker development. *Cancer*
727 *Discov.*, **3**,148–157.
- 728 46. Leong, T.Y.-M. and Leong, A.S.-Y. (2005) Epidemiology and carcinogenesis of
729 hepatocellular carcinoma. *HPB (Oxford)*, **7**,5–15.
- 730 47. Rothwell, J., Fillâtre, Y., Martin, J.-F., Lyan, B., Pujos-Guillot, E., Fezeu, L., Hercberg,
731 S., Comte, B., Galan, P., Touvier, M., and Manach, C. (2014) New biomarkers of
732 coffee consumption identified by the non-targeted metabolomic profiling of cohort
733 study subjects. *PLoS One*, **9**,e93474.
- 734 48. Guo, M., Zhao, B., Liu, H., Zhang, L., Peng, L., Qin, L., Zhang, Z., Li, J., Cai, C., and
735 Gao, X. (2014) A Metabolomic Strategy to Screen the Prototype Components and
736 Metabolites of Shuang-Huang-Lian Injection in Human Serum by Ultra Performance
737 Liquid Chromatography Coupled with Quadrupole Time-of-Flight Mass Spectrometry.
738 *J. Anal. Methods Chem.*, **2014**,241505.
- 739 49. Anderson, T.W. (1951) Estimating linear restrictions on regression coefficients for
740 multivariate normal distributions. *Ann. Math. Stat.*, **22**,327–351.
- 741 50. Izenman, A.J. (1975) Reduced-Rank Regression for the Multivariate Linear Model. *J.*
742 *Multivar. Anal.*, **5**,248–264.
- 743 51. Aldrin, M. (2002) Reduced-rank regression. In: El-Shaarawi AH, Piegorsch WW,
744 editors. Encyclopedia of Environmetrics. Chichester: John Wiley & Sons, Ltd. pp.
745 1724–1728.
- 746 52. Hotelling, H. (1936) Relations between two sets of variates. *Biometrika*, **28**,321–377.
- 747 53. Navratil, V., Pontoizeau, C., Billoir, E., and Blaise, B.J. (2013) SRV : an opensource
748 toolbox to accelerate the recovery of metabolic biomarkers and correlations from
749 metabolic phenotyping data sets. *Bioinformatics*, **29**,1348–1349.
- 750 54. Wishart, D.S. (2008) Quantitative metabolomics using NMR. *TrAC Trends Anal.*
751 *Chem.*, **27**,228–237.
- 752 55. Psychogios, N., Hau, D.D., Peng, J., Guo, A.C., Mandal, R., Bouatra, S., Sinelnikov, I.,
753 Krishnamurthy, R., Eisner, R., Gautam, B., Young, N., Xia, J., Knox, C., Dong, E.,

754 Huang, P., Hollander, Z., Pedersen, T.L., Smith, S.R., Bamforth, F., Greiner, R.,
755 McManus, B., Newman, J.W., Goodfriend, T., and Wishart, D.S. (2011) The human
756 serum metabolome. *PLoS One*, **6**,e16957.

757

758 **Legends to figures**

759 **Figure 1:** General scheme of the analytical framework developed in the study. A PC-PR2
760 analysis is carried out beforehand to identify relevant sources of variation. In the PLS model
761 the X- and Y- sets are related to each other, and scores are computed (1). X- and Y-scores are,
762 in turn, associated to a case-control indicator of HCC status in conditional logistic regression
763 models (2). A mediation analysis is carried out to explore the role of metabolomics in the
764 association between lifestyle factors and risk of HCC (3).

765 **Figure 2:** PC-PR2 analysis results* identifying the sources of variability in the NMR data
766 (panel A) and in the lifestyle data (panel B).

767 * 17 and 14 components were retained to account for 80 % (threshold used) of total NMR (A)
768 and lifestyle variability (B), respectively. The R2 value represents the amount of variability in
769 NMR / lifestyle variable explained by the ensemble of investigated predictors.

770 **Figure 3:** Variable importance plot (VIP) displaying the variable importance for projection
771 statistic of the predictor variables for the PLS analyses.

772 Panel A: Results from the main PLS model run on all observations (N=336, X-set=21, Y-
773 set=285).

774 Panel B: Results from the PLS sensitivity analysis run on a subsample (N=271, 92 cases, 179
775 controls) excluding sets where cases were diagnosed within the first two years of follow-up
776 (X-set=21, Y-set=285).

777 The horizontal line corresponds to Wold's criterion (0.8), the threshold used to rule if a
778 variable has an important contribution to the construction of the Y variables (see
779 **Mathematical Appendix** for further details).

Table 1: Individual and cumulative variation (%) explained by the first 3 PLS factors in 21 lifestyle (X-set) and 285 NMR (Y-set) variables.

# of PLS Factors	Lifestyle Variables		NMR Variables	
	Individual	Cumulative	Individual	Cumulative
1	6.17	-	5.51	-
2	6.23	12.40	2.38	7.89
3	9.27	21.67	0.59	8.48

Table 2: Lifestyle and NMR cluster variables contributing to each of the 3 PLS factors (N=336, X-set=21, Y-set=285).

PLS Factor	Lifestyle Variable*	Loading value	CS* \ddagger (ppm)	Metabolite**	Loading value
1	Dairy Products	0.28	5.22		-0.06
	Cakes and Biscuits	0.32	3.88		-0.05
	Lifetime Alcohol Consumption	-0.25	3.82		-0.06
	Smoking Status	-0.39	3.76		-0.06
	Diabetes	-0.63	3.71	Glucose	-0.05
			3.54		-0.05
			3.50		-0.07
			3.48		-0.07
			3.44	Acetoacetate	-0.07
			3.23	Choline + Glycerphosphocholine	-0.04
			3.01	Lysine	0.10
			2.94	Albumin	0.10
			2.65	Aspartate	0.10
			2.42	Glutamine	0.10
			2.28	Acetoacetate	0.10
			2.22	CH ₂ -CH ₂ -COOC bond of lipids + Acetone	-0.04
			1.86		0.09
			1.87	Lysine	0.10
			1.53	CH ₂ -CH ₂ -COOC bond of lipids	-0.03
2	Cereal and Cereal Products	-0.16	7.17		0.13
	Height	-0.34	6.87	Tyrosine	0.13
	Education Level	-0.26	5.27	CH=CH bond of lipids	-0.13
	Hepatitis	0.49	5.22	Glucose	0.16
			5.18	Mannose + Lipid O-CH ₂	-0.12
			4.27	Lipid O-CH ₂	-0.12
			4.25	Threonine	-0.14
			4.07	Choline + Lipid O-CH ₂ + Myo-inositol	-0.12
			4.05	Creatinine	-0.14
			3.88		0.15
			3.82		0.16
			3.76		0.15
			3.71	Glucose	0.15
			3.54		0.15
			3.50		0.16
			3.48		0.16
			3.44	Acetoacetate	0.16
			3.23	Choline + Glycerphosphocholine	0.15
			2.80	Aspartate	-0.12
			2.22	CH ₂ -CH ₂ -COOC bond of lipids + Acetone	-0.11
			2.19	CH ₂ -CH ₂ -COOC bond of lipids	-0.15
		2.02	Proline + Glutamate + CH ₂ =C bonds of lipids	-0.13	
		1.53	CH ₂ -CH ₂ -COOC bond of lipids	-0.13	
		1.25	CH ₂ bond of lipids	-0.12	
		0.86	Cholesterol + CH ₃ bond of lipids	-0.12	
3	Vegetables	-0.42	7.32	Phenylalanine	0.11
	Lifetime Alcohol Consumption	0.29	5.22	Glucose	-0.13
	Smoking Status	0.25	4.28	Lipid O-CH ₂	0.11
	Hepatitis	0.26	3.88		-0.11
			3.82		-0.11
			3.76	Glucose	-0.12
			3.71		-0.11
			3.69		-0.11
			3.63	Myo-inositol	0.16
			3.50	Glucose	-0.13
			3.48		-0.12
			3.44	Acetoacetate	-0.12
			3.35	Proline	0.11
			3.33		0.13
			3.28	Myo-inositol	0.12
			3.23	Choline + Glycerphosphocholine	-0.12
			2.80	Aspartate	-0.13
			2.76	part of =CH-CH ₂ -CH= bond of lipids	-0.13
			2.35		0.12
			2.33	Proline + Glutamate	0.13

	1.20	3-hydroxybutyrate + CH2 bond of lipids	0.11
	1.16	Ethanol	0.15
	0.66	Cholesterol	0.11

*Relevant lifestyle and NMR variables contributing to each PLS factor selected based on their associated loading values <10th percentile (pctl) and >90th pctl or <5th pctl and >95th pctl respectively.

‡ CS: ¹H chemical shift (in ppm) of the cluster (center value).

**Some of the identified clusters were found to be background noise during the annotation phase and were removed from this table.

Table 3: HCC odds ratios* and 95% confidence interval (OR, 95% CI) associated with the lifestyle (X-set) and the NMR clusters (Y-set) PLS scores in the main analysis (N=336, X-set=21, Y-set=285).

PLS Lifestyle Variables X-scores			PLS NMR Variables Y-scores		
Factor	OR** (95% CI)	P-Wald†	Factor	OR** (95% CI)	P-Wald†
1	0.77 (0.58, 1.02)	0.07	1	0.96 (0.91, 1.01)	0.09
2	1.54 (1.06, 2.25)	0.02	2	1.11 (1.02, 1.22)	0.02
3	1.37 (1.05, 1.79)	0.02	3	1.22 (1.04, 1.44)	0.01

*Models were adjusted for C-reactive protein concentration, alpha-fetoprotein concentration and a composite score for liver damage. Cases and controls were matched on age at blood collection (± 1 year), sex, study centre, date (± 2 months) and time of the day at blood collection (± 3 hours), fasting status at blood collection (<3/3-6/>6 hours); among women, additional matching criteria included menopausal status (pre-/peri-/postmenopausal) and hormone replacement therapy use at time of blood collection (yes/no). ** ORs expressing the change in HCC risk associated to 1-SD increase in the score. † Wald's test was for continuous exposure compared with a Chi-square distribution with 1 degree of freedom (dof).

Table 4: Area under the curve (AUC), sensitivity, specificity, accuracy and positive predictive value (PPV) of ROC models (with 95% CI), from the main PLS analysis (N=336, X-set=21, Y-set=285).

	AUC	AUC _b **	Sensitivity	Specificity	Accuracy	PPV
Adjustment Covariates (ADJ)*	0.842 (0.794, 0.891)	0.821 (0.766, 0.868)	0.752 (0.662, 0.829)	0.802 (0.743, 0.852)	0.785	0.0015
X1 scores + ADJ	0.846 (0.797, 0.894)	0.825 (0.766, 0.875)	0.743 (0.653, 0.821)	0.838 (0.783, 0.884)	0.806	0.0018
X1+X2 scores + ADJ	0.854 (0.808, 0.900)	0.831 (0.772, 0.881)	0.743 (0.653, 0.821)	0.824 (0.768, 0.872)	0.797	0.0017
X1+X2+X3 scores + ADJ	0.859 (0.811, 0.907)	0.836 (0.778, 0.887)	0.796 (0.710, 0.866)	0.788 (0.729, 0.840)	0.791	0.0015
Y1 scores + ADJ	0.841 (0.793, 0.890)	0.817 (0.760, 0.865)	0.735 (0.643, 0.813)	0.820 (0.763, 0.868)	0.791	0.0016
Y1+Y2 scores + ADJ	0.845 (0.795, 0.894)	0.820 (0.762, 0.872)	0.735 (0.643, 0.813)	0.851 (0.798, 0.895)	0.812	0.0020
Y1+Y2+Y3 scores + ADJ	0.853 (0.804, 0.902)	0.827 (0.771, 0.877)	0.726 (0.634, 0.805)	0.883 (0.833, 0.922)	0.890	0.0025

*The model is run on the adjustment covariates (ADJ) including the C-reactive protein concentration, alpha-fetoprotein concentration and a composite score for liver damage. ** AUC_b is the bootstrapped-cross validated estimate of the AUC. X1, X2 and X3 are the lifestyle component scores of the first, second and third PLS factors, respectively. Y1, Y2, and Y3 are the metabolomics component of the first, second and third PLS factors, respectively.

Table 5: HCC odds ratios* and 95% confidence intervals (OR, 95%CI) associated with the lifestyle (X-set) and the NMR clusters (Y-set) PLS scores. Results from the sensitivity analysis (N=271, 92 cases, 179 controls) conducted excluding sets where cases were diagnosed within the first two years of follow-up (X-set=21, Y-set=285).

PLS Lifestyle Variables X-scores			PLS NMR Variables Y-scores		
Factor	OR** (95% CI)	P-Wald†	Factor	OR** (95% CI)	P-Wald†
1	0.80 (0.60, 1.08)	0.15	1	0.96 (0.94, 1.04)	0.56
2	1.56 (1.02, 2.40)	0.04	2	1.18 (1.03, 1.36)	0.02
3	0.86 (0.67, 1.11)	0.26	3	0.86 (0.73, 0.99)	<0.05

*Models were adjusted for C-reactive protein concentration, alpha-fetoprotein concentration and a composite score for liver damage. Cases and controls were matched on age at blood collection (± 1 year), sex, study centre, date (± 2 months) and time of the day at blood collection (± 3 hours), fasting status at blood collection (<3/3-6/>6 hours); among women, additional matching criteria included menopausal status (pre-/peri-/postmenopausal) and hormone replacement therapy use at time of blood collection (yes/no). ** ORs expressing the change in HCC risk associated to 1-SD increase in the score. † Wald's test was for continuous exposure compared with a Chi-square distribution with 1 degree of freedom (dof).

Table 6: Results from the mediation analysis (N= 336, X-set=21, Y-set=285): Natural Indirect Effect (NIE) and 95%CI*.

Model**				Natural Indirect Effect (NIE)	
Exposure (A)	Mediator (M)	Outcome	A*M interaction term	Estimate (95%CI)	p-value
X1 score	Y1 score	HCC	No	0.91 (0.77, 1.06)	0.23
X2 score	Y2 score	HCC	No	1.11 (0.97, 1.25)	0.12
X3 score	Y3 score	HCC	No	1.08 (0.94, 1.23)	0.28
X1 score	Y1 score	HCC	Yes	0.96 (0.79, 1.17)	0.70
X2 score	Y2 score	HCC	Yes	1.15 (1.01, 1.31)	0.04
X3 score	Y3 score	HCC	Yes	1.13 (1.01, 1.28)	0.04

* The standard errors used to compute the 95%CI were obtained using the delta method.

**Models were adjusted for the C-reactive protein concentration, alpha-fetoprotein concentration and a composite score for liver damage, as well as for the other Y-scores, as potential mediator-outcome confounders. Additionally, the outcome and the mediator models were adjusted for centre and age at blood collection.

Mathematical Appendix

A statistical framework to model the meeting-in-the-middle principle using metabolomic data: application to hepatocellular carcinoma in the EPIC study.

1 PLS regression

1.1 Introduction

PLS (partial least squares) regression is a widely used method in multivariate statistics to relate two sets of variables while reducing their dimensionality. It was first developed as a method to predict a set of variables Y from another set X ; and also to depict their common structure. The main aim of PLS is to regress a set Y of \mathbf{q} variables (y_1, y_2, \dots, y_q) of interest, which are called responses, on a set X of \mathbf{p} predictor variables (x_1, x_2, \dots, x_p) that may display high levels of correlation. PLS combines and generalizes features of principal component analysis (PCA) and multiple linear regression (MLR); and results in a set of PLS latent factors as linear combinations of variables, in turn, in the X - and Y -sets. By simultaneously decomposing X and Y , PLS finds components that explain as much as possible of the inter-relations of X and Y . The latent factors obtained from the decomposition can be used to predict Y . The following details of the algorithm are adapted from Michel Tenenhaus' book *La régression PLS, Théorie et Pratique* [1].

1.2 The PLS algorithm

Two different, but closely related, techniques exist under the name of PLS regression. The canonical or symmetric PLS regression assumes that the X - and Y - sets play a symmetrical role. The version presented here is the regression mode where latent variables are computed from a succession of singular value decompositions (SVD) followed by deflation of both the X - and Y - matrices. These sets are assumed to play the asymmetric roles of predictors and responses, respectively. Next, we briefly describe the landmark algorithm NIPALS Nonlinear estimation by Iterative Partial Least Squares. As a first step, two substitute matrices X_0 and Y_0 are initialized with $X_0 = X_{(n \times p)}$ and $Y_0 = Y_{(n \times q)}$, where variables were standardized to have means and standard deviations equal to zero and one, respectively. For $h = 1, \dots, H$, where $H = \min(p, q)$, the PLS factors are obtained iteratively. PLS regression focuses on finding two sets of weights, $w_{h(p \times 1)}$ and $c_{h(q \times 1)}$, in order to create respectively a linear combination of the columns of X and Y , known as the PLS factors, such that these two linear combinations have maximum covariance and are unique. These weights define a first pair of vectors, called the X - and Y -scores, $t_h = Xw_h$ and $u_h = Yc_h$ where we have $t_h^\top u_h$ maximal. PLS can be written as the following optimisation problem where maximum covariance is sought between $t_{h(1 \times n)}$ and $u_{h(1 \times n)}$ for each $h = 1 \dots H$:

$$\text{Max } cov(Xw_h, Yc_h) \quad (1)$$

under the following normality constraints

$$\|w_h\| = 1 \quad (2)$$

$$\|c_h\| = 1 \quad (3)$$

and the following orthogonality constraint

$$t_h^\top(t_1, \dots, t_{h-1}) = 0 \quad (4)$$

By construction we also have the following property:

$$u_h^\top(t_1, \dots, t_{h-1}) = 0 \quad (5)$$

The first pair of X - and Y - scores can equivalently be obtained via a singular value decomposition. Indeed, the SVD of the cross-product matrix $X_{h-1}^\top Y_{h-1}$ leads to the identification of the first left and right singular vectors and of the weights w_h and c_h . The scores t_h and u_h are obtained as follows:

$$t_h = X_{h-1} w_h \quad (6)$$

$$u_h = Y_{h-1} c_h \quad (7)$$

The vector t_h is then normalized (a scaling of u_h is optional). Regressing the predictor and response matrices on the t_h vector yields the corresponding loadings.

$$p_h = X_{h-1}^\top t_h \quad (8)$$

$$c_h = Y_{h-1}^\top t_h \quad (9)$$

Next is the deflation step, where information based on the extracted latent factor h is subtracted from the current data matrices.

$$X_h = X_{h-1} - t_h p_h^\top \quad (10)$$

$$Y_h = Y_{h-1} - t_h c_h^\top \quad (11)$$

The described steps of the algorithm are iterated until one of the following criteria is met:

- If H is specified, and the algorithm stops when the H -th PLS factor is extracted and its associated statistics computed.
- If H is not specified, the algorithm stops when X_H becomes a null matrix. In this case however, H cannot exceed $\min(p, q)$.

Algorithm 1 PLS1 classic algorithm steps - When Y is univariate.

- 1: $X_0 \leftarrow X$; $y_0 \leftarrow y$
 - 2: **for** ($h = 1$; $h \leq H$; $h++$) **do**
 - 3: $w_h = X_{h-1}^\top y_{h-1} / y_{h-1}^\top y_{h-1}$
 - 4: $w_h = w_h / \sqrt{w_h^\top w_h}$
 - 5: $t_h = X_{h-1} w_h / w_h^\top w_h$
 - 6: $p_h = X_{h-1}^\top t_h / t_h^\top t_h$
 - 7: $X_h = X_{h-1} - t_h p_h^\top$
 - 8: $c_h = y_{h-1}^\top t_h / t_h^\top t_h$
 - 9: $u_h = y_{h-1} / c_h$
 - 10: $y_h = y_{h-1} - c_h t_h$
-

When Y is univariate, the PLS algorithm carried out is PLS1 (See Algorithm 1, following the notation of M. Tenenhaus [1]). PLS2 (Algorithm 2) is used when Y is multivariate. When there are missing data in either the X - or Y - sets, the coordinates of the vectors w_h , t_h , c_h , u_h , and p_h are computed as slopes of the least squares straight line that passes through the origin, using the available data as follows:

Algorithm 2 PLS2 classic algorithm steps - When Y is multivariate.

```
1:  $X_0 \leftarrow X$  ;  $Y_0 \leftarrow Y$ 
2: for ( $h = 1$ ;  $h \leq H$ ;  $h++$ ) do
3:    $u_h = Y_{h-1}[, 1]$  i.e. the first column of the matrix
4:   while  $w_h$  has not converged do
5:      $w_h = X_{h-1}^\top u_h / u_h^\top u_h$ 
6:      $w_h = w_h / \sqrt{w_h^\top w_h}$ 
7:      $t_h = X_{h-1} w_h / w_h^\top w_h$ 
8:      $c_h = Y_{h-1}^\top t_h / t_h^\top t_h$ 
9:      $u_h = Y_{h-1} c_h / c_h^\top c_h$ 
10:     $p_h = X_{h-1}^\top t_h / t_h^\top t_h$ 
11:     $X_h = X_{h-1} - t_h p_h^\top$ 
12:     $Y_h = Y_{h-1} - t_h c_h^\top$ 
```

- $w_h = (w_{h1}, \dots, w_{hp})^\top$, is a normalized vector, where w_{hj} is the slope of the least squares line passing through the origin of the plane defined by $(u_h, X_{h-1,j})$. $X_{h-1,j}$ is the j -th X variable of the $h - 1$ PLS factor.
- $t_h = (t_{h1}, \dots, t_{hn})^\top$, where t_{hi} is the slope of the least squares line passing through the origin of the plane defined by $(w_h, x_{h-1,i})$. $x_{h-1,i}$ is the i -th x observation of the $h - 1$ PLS factor.
- $c_h = (c_{h1}, \dots, c_{hq})^\top$, where c_{hk} is the slope of the least squares line passing through the origin of the plane defined by $(t_h, Y_{h-1,k})$. $Y_{h-1,k}$ is the k -th Y variable of the $h - 1$ PLS factor.
- $u_h = (u_{h1}, \dots, u_{hn})^\top$, where u_{hi} is the slope of the least squares line passing through the origin of the plane defined by $(c_h, y_{h-1,i})$. $y_{h-1,i}$ is the i -th y observation of the $h - 1$ PLS factor.
- $p_h = (p_{h1}, \dots, p_{hp})^\top$, where p_{hj} is the slope of the least squares line passing through the origin of the plane defined by $(t_h, X_{h-1,j})$. $X_{h-1,j}$ is the j -th X variable of the $h - 1$ PLS factor.

1.3 Tools for interpretation

1.3.1 Choice of number of components

The number of PLS latent factors or components to be retained can be decided based on a cross-validation.

For each model with a number h of extracted factors, this is done by running the PLS analysis on only a part of the data called the training set, and then evaluating how well the model fits observations in the test set. This includes the part of the data not involved in the PLS modelling of the training set.

The dataset comprised of n observations is split into z approximately equal sets of observations. The training set consists of the data in the first $z - 1$ folds and the remaining fold is used as test set. Predicted values for the Y -set are computed on this test set along with the sum of the squared error of prediction. This process is repeated z times so that each fold can in turn serve as a test set. In practice, for each number of possible latent factors $h = 1, \dots, H$, we compute the prediction of y_i by the PLS model with results obtained on the training set with a number h of components applied to observations in the test set in order to yield $\hat{y}_{h(-i)}$. The Prediction Error Sum of Squares (PRESS) is the resulting sum of all squared errors of prediction statistic computed across all test sets as defined in the following equation:

$$PRESS_h = \sum (y_i - \hat{y}_{h(-i)})^2 \quad (12)$$

The Residual Sum of Squares (RSS) is computed in a standard way:

$$RSS_h = \sum (y_i - \hat{y}_{hi})^2 \quad (13)$$

Different criteria can be used to determine the number of components h to retain. One such criterion, Q_h^2 was first introduced by H. Wold [2] and is mainly used in the software SIMCA-P. It is based on the following statistic:

$$Q_h^2 = 1 - \frac{PRESS_h}{RSS_{h-1}} \quad (14)$$

As pointed out by M. Tenenhaus, the initial value for RSS when y is univariate centred-scaled and $h = 0$ is:

$$RSS_0 = \sum_{i=1}^n (y_i - \bar{y})^2 = n - 1 \quad (15)$$

In the software SIMCA-P the PLS component is kept when the following condition is met:

$$\sqrt{PRESS_h} \leq 0.95\sqrt{RSS_{h-1}} \quad (16)$$

$$\iff Q_h^2 \geq 0.0975 \quad (17)$$

The default threshold 0.0975 is equal to $1 - 0.95^2$. In SAS, the criteria to select the number h of components to be retained is by minimizing the $PRESS_h$ statistic.

The above described formulae can be generalized for multivariate Y , thus we have for any given variable y_k , $k = 1, \dots, q$:

$$Q_{kh}^2 = 1 - \frac{PRESS_{kh}}{RSS_{k(h-1)}} \quad (18)$$

$$Q_h^2 = 1 - \frac{\sum_{k=1}^q PRESS_{kh}}{\sum_{k=1}^q RSS_{k(h-1)}} \quad (19)$$

The criteria for keeping a PLS factor are identical to what was established for the univariate case. One can alternately use one of the following rules, where the equivalence defined in formula (17) still holds true:

- $Q_h^2 \geq 0.0975$
- At least one value of $Q_{hk}^2 \geq 0.0975$

If the criteria are met by several values of h , the one retained is the smallest h , to achieve a better dimensionality reduction.

The Q^2 and $PRESS$ criteria are relatively robust to the choice of number of folds (blocks) used for cross-validation. A number of folds between 5 and 10 is recommended (Tenenhaus 1998, p.238) [1]. The default choice in the SIMCA-P and SAS softwares is 7, and is the parameter used in this study.

1.3.2 Variable Importance in the Projection (VIP)

The Variable Importance in the Projection (VIP) is a measure of the explanatory power of a given variable x_j over Y . The VIP_{hj} of a given component h of the j -th variable x_j is defined as:

$$VIP_{hj} = \sqrt{\frac{p}{Rd(Y; t_1, \dots, t_h)} \sum_{l=1}^h Rd(Y, t_l) w_{lj}^2} \quad (20)$$

and one has:

$$\sum_{j=1}^p VIP_{hj}^2 = p \quad (21)$$

where $Rd(Y; t_1, \dots, t_h)$ is the redundancy of Y with respect to the t scores (t_1, \dots, t_h) . It describes the amount of variance of Y explained by the component t_h of the X -set. It is defined

as follows:

$$Rd(Y, t_h) = \frac{1}{q} \sum_{k=1}^q cor^2(y_k, t_h) \quad (22)$$

It can be equivalently computed as:

$$Rd(Y, t_h) = r_h^2 \frac{1}{q} \sum_{k=1}^q cor^2(y_k, u_h) \quad (23)$$

where $r_h = cor(Xw_h, Yc_h)$ is called a canonical correlation and r_h^2 is the h^{th} largest eigenvalue of the crossproduct matrix decomposition.

The contribution of a variable x_j to the construction of a component t_l is measured by the weight w_{lj}^2 . For each l , with $l = 1, \dots, h$, the sum of these weights across the p variables x_j equals 1. To measure the contribution of the variable x_j to the construction of Y through the components t_l , one should consider the explanatory power of the component t_l , measured by the redundancy $Rd(Y; t_l)$. An equal weight w_{lj}^2 indicates an explanatory power of the x_j variable over the Y -set whose importance increases with the level of redundancy $Rd(Y; t_l)$.

The VIP enables the ranking of the predictors x_j according to their explanatory power on Y , and summarizes their contribution to the model. A VIP is considered small if its value is less than 0.8 and high when its value is greater than 1. Variables with a high VIP ($VIP > 1$) are the most important for the reconstruction and prediction of Y .

2 Statistical Recoupling of Variables (SRV)

The SRV procedure was introduced by *Blaise et al. (2009)* [3] and for which a matlab toolbox was later implemented [4]. The SRV is an "intelligent bucketing" algorithm that aims at regrouping variables (typically the smallest unit of the NMR spectrum) in clusters corresponding to a wider biological and chemical entity.

SRV exploits the spectral structure of data, without forming any metabolic hypothesis to reduce the dimensionality of spectra. A typical NMR 1H 9 ppm spectrum is often partitioned into 9,000 buckets of 0.001 ppm width. The main idea of the algorithm is to exploit the spectral dependency landscape L which is the covariance to correlation ratio between two neighbouring variables along the chemical shift axis to assemble them within a cluster. If one considers a matrix Z of serum spectra acquired by NMR with n observations and r columns (z_1, \dots, z_r) corresponding to neighbouring bins of NMR signals. The first bin-variable starts the first

cluster, then L is computed for each z_i as follows with $i = 1, \dots, r$:

$$\begin{aligned} L(z_i) &= \frac{\text{cov}(z_i, z_{i+1})}{\text{cor}(z_i, z_{i+1})} \\ &= \text{sd}(z_i) * \text{sd}(z_{i+1}) \end{aligned} \tag{24}$$

where sd is the standard deviation.

The variable then joins a cluster according to the following rules:

- $L(z_i)$ values are used to locate local minima i.e. borders between clusters.
- If $L(z_{i-1}) > L(z_i)$ then z_{i-1} and z_i are associated in the same cluster, otherwise z_i and z_{i+1} start a new cluster.
- The minimum number of variables belonging to a cluster is set a priori as it is based on the resolution of the NMR spectra. When acquired at 700 MHz, the typical peak base width of a well-resolved singlet is equal to 7 Hz. Therefore, the threshold was set to 10 in our analysis, meaning that if a cluster has less than 10 variables, it is discarded.
- The super-cluster intensity is computed as the mean of the intensities of the signal in the bins assigned to the super-cluster.
- If two neighbouring clusters have a correlation > 0.9 , they are aggregated to form a super-cluster. In these analyses, the association is limited to 3 clusters per super-cluster (this value is empirical and was discussed in the original paper [3]).

References

- [1] Michel Tenenhaus. *La régression PLS: Théorie et Pratique*. Paris, 1998.
- [2] Herman Wold. *Multivariate Analysis*, chapter Estimation of principal components and related models by iterative least squares, pages 391–420. Academic Press, New York, 1966.
- [3] Benjamin J Blaise, Laetitia Shintu, Bénédicte Elena, Lyndon Emsley, Marc-Emmanuel Dumas, and Pierre Toulhoat. Statistical recoupling prior to significance testing in nuclear magnetic resonance based metabonomics. *Analytical Chemistry*, 81(15):6242–6251, August 2009.
- [4] Vincent Navratil, Clément Pontoizeau, Elise Billoir, and Benjamin J Blaise. Srv: an open-source toolbox to accelerate the recovery of metabolic biomarkers and correlations from metabolic phenotyping data sets. *Bioinformatics*, 29(10):1348–1349, May 2013.

1 **A Statistical framework to model the meeting-in-the-middle principle using**
2 **metabolomic data: application to hepatocellular carcinoma in the EPIC study.**

3 *Nada Assi¹, Anne Fages^{2,a}, Paolo Vineis³, Marc Chadeau-Hyam³, Magdalena Stepień¹, Talita*
4 *Duarte-Salles¹, Graham Byrnes¹, Houda Boumaza², Sven Knüppel⁴, Tilman Kühn⁵, Domenico*
5 *Palli⁶, Christina Bamia⁷, Hendriek Boshuizen⁸, Catalina Bonet⁹, Kim Overvad¹⁰, Mattias*
6 *Johansson^{1,11}, Ruth Travis¹², Marc J Gunter³, Eiliv Lund¹³, Laure Dossus¹⁴⁻¹⁵, Bénédicte*
7 *Elena-Herrmann², Elio Riboli³, Mazda Jenab¹, Vivian Viallon^{16-18†}, Pietro Ferrari^{1†*}*

8

9 ¹ International Agency for Research in Cancer (IARC-WHO), 150 Cours Albert Thomas,
10 69372 Lyon CEDEX 08, France

11 ² Centre de RMN à Très Hauts Champs, Institut des Sciences Analytiques (CNRS/ENS
12 Lyon/UCB Lyon 1), Université de Lyon, 69100 Villeurbanne, France

13 ³ Department of Epidemiology and Biostatistics, MRC-HPA Centre for Environment and
14 Health, School of Public Health, Imperial College London, Norfolk Place, London, W2 1PG,
15 United Kingdom

16 ⁴ Department of Epidemiology, German Institute of Human Nutrition, Potsdam-Rehbrücke,
17 14558 Nuthetal, Germany

18 ⁵ Division of Cancer Epidemiology, German Cancer Research Center (DKFZ), Heidelberg,
19 Germany

20 ⁶ Molecular and Nutritional Epidemiology Unit, Cancer Research and Prevention Institute –
21 ISPO, Florence- Italy

22 ⁷ WHO Collaborating Center for Food and Nutrition Policies, Department of Hygiene,
23 Epidemiology and Medical Statistics, University of Athens Medical School, Athens, Greece

24 ⁸ National Institute for Public Health and the Environment (RIVM), Bilthoven, The
25 Netherlands

26 ⁹ Unit of Nutrition and Cancer, Cancer Epidemiology Research Program, Institut Català
27 d'Oncologia, L'Hospitalet de Llobregat, Spain

28 ¹⁰ The Department of Epidemiology, School of Public Health, Aarhus University, Aarhus,
29 Denmark

30 ¹¹ The Department for Biobank Research, Umeå University, Sweden

31 ¹² Cancer Epidemiology Unit, Nuffield Department of Population Health University of
32 Oxford, Oxford, UK

33 ¹³ The Institute of Community Medicine, University of Tromsø, Tromsø, Norway

34 ¹⁴ Inserm, Centre for research in Epidemiology and Population Health (CESP), U1018,
35 Lifestyle, genes and health: integrative trans-generational epidemiology team, Villejuif,
36 France

37 ¹⁵ Université Paris Sud, Villejuif, France

38 ¹⁶ Université de Lyon, F-69622, Lyon, France

39 ¹⁷ Université Lyon 1, UMRESTTE, F-69373 Lyon, France

40 ¹⁸ IFSTTAR, UMRESTTE, F-69675 Bron, France

41 ^a Present address: Chemical Physics Department, Weizmann Institute of Science, Rehovot,
42 Israel

43 [†] Joint last authors.

44

45 ^{*}Corresponding author: Pietro Ferrari, International Agency for Research on Cancer, 150
46 Cours Albert Thomas, 69372 Lyon cedex 08, France. Tel: +33 472 73 8031; Fax: +33 472 73
47 8361. E-mail: ferrarip@iarc.fr

48 **Abstract**

49 Metabolomics is a potentially powerful tool for identification of biomarkers associated with
50 lifestyle exposures and risk of various diseases. This is the rationale of the “meeting-in-the-
51 middle” concept, for which an analytical framework was developed in this study. In a nested
52 case-control study on hepatocellular carcinoma (HCC) within the European Prospective
53 Investigation into Cancer and nutrition (EPIC), serum ¹H NMR spectra (800 MHz) were
54 acquired for 114 cases and 222 matched controls. Through Partial Least Square (PLS)
55 analysis, 21 lifestyle variables (the “predictors”, including information on diet, anthropometry
56 and clinical characteristics) were linked to a set of 285 metabolic variables (the “responses”).
57 The three resulting scores were related to HCC risk by means of conditional logistic
58 regressions. The first PLS factor was not associated with HCC risk. The second PLS
59 metabolomic factor was positively associated with tyrosine and glucose, and was related to a
60 significantly increased HCC risk with OR= 1.11 (95%CI: 1.02, 1.22, p=0.02) for a 1-SD
61 change in the responses score, and a similar association was found for the corresponding
62 lifestyle component of the factor. The third PLS lifestyle factor was associated with lifetime
63 alcohol consumption, hepatitis and smoking, and had negative loadings on vegetables intake.
64 Its metabolomic counterpart displayed positive loadings on ethanol, glutamate and
65 phenylalanine. These factors were positively and statistically significantly associated with
66 HCC risk, with 1.37 (1.05, 1.79, p=0.02) and 1.22 (1.04, 1.44, p=0.01), respectively. Evidence
67 of mediation was found in both the second and third PLS factors, where the metabolomic
68 signals mediated the relation between the lifestyle component and HCC outcome. This study
69 devised a way to bridge lifestyle variables to HCC risk through NMR metabolomics data.
70 This implementation of the “meeting-in-the-middle” approach finds natural applications in
71 settings characterized by high-dimensional data, increasingly frequent in the –omics
72 generation.

- 73 **Keyword:** partial least square, lifestyle factors, metabolomics, hepatocellular carcinoma,
74 meeting-in-the-middle, molecular epidemiology.

75 **Introduction**

76 Metabolomic profiles from blood and other biological samples collected from large-
77 scale epidemiologic studies are increasingly being investigated [1], following recent
78 developments in nuclear magnetic resonance (NMR) and mass spectrometry (MS) enabling
79 the assessment of metabolic profiles for large numbers of individuals. As a result,
80 metabolomic data is gradually playing a key part in clinical and observational studies; and
81 new statistical methodologies [2] are increasingly being sought to explore insights into
82 pathological processes that metabolomics may provide in order to better understand
83 determinants of disease development. These approaches explore a variety of etiological
84 hypotheses; however they usually focus on one aspect at a time, combining metabolomics
85 with either epidemiologic/phenotypic data on lifestyle exposures [3] or with disease outcomes
86 [4,5]. The main aim of this work is to jointly use all aspects that are potentially informative to
87 apprehend the contrivances of disease development.

88 Metabolomic data offers the opportunity to identify signatures and biomarkers
89 associated with environmental exposures and the risk of a disease. Prospective studies are
90 conceptually suitable for this purpose, since they rely on biological samples collected before
91 disease onset, and are thus marginally influenced by metabolic changes due to processes of
92 disease development. In this scenario, the “meeting-in-the-middle” (MITM) approach [6] has
93 been conceived as a research strategy to identify biomarkers that are related to specific
94 exposures and that are, at the same time, predictive of disease outcome. Finding this overlap
95 between exposure and disease of “intermediate” biomarkers can potentially disclose useful
96 information on the exposure-to-disease pathway, and may serve as an objective risk exposure
97 measure, ultimately allowing the identification of a targeted prevention scheme. The MITM
98 was previously implemented as a proof of concept in a case-control study nested within a
99 cohort of healthy individuals [7], where a list of putative intermediate ¹H NMR biomarkers

100 linking exposure to dietary compounds, mainly micro- and macronutrients, and disease
101 outcomes (colon and breast cancer) were investigated.

102 In this study we extend previous attempts to model the MITM by fully integrating
103 metabolomics, lifestyle and disease risk in a single analytical framework. A strategy was
104 developed to simultaneously investigate a broad range of metabolites and lifestyle variables
105 with a partial least square (PLS) regression model [8]. The resulting scores were related to the
106 risk of hepatocellular carcinoma (HCC), in a case-control study nested within the European
107 Prospective Investigation into Cancer and nutrition (EPIC). HCC is the most frequent primary
108 form of cancer affecting the liver, an organ that plays a critical role in many metabolic
109 pathways [9]. HCC is a disease with multifactorial origins embracing lifestyle and dietary
110 exposures whose intersection may reveal metabolomic signals [10] relevant to cancer onset.
111 The system of relationships between metabolomic profiles and lifestyle factors in relation to
112 HCC was evaluated by means of mediation analysis. The methodological challenges
113 characterizing the analysis of large and complex metabolomic datasets are described and
114 discussed.

115 **Methods**

116 *EPIC design.* The European Prospective Investigation into Cancer and nutrition (EPIC) is a
117 large cohort established to investigate the association of diet, lifestyle and environmental
118 factors with cancer incidence and other chronic disease outcomes. Between 1992-2000, over
119 520,000 participants aged 20-85 years, were recruited from 23 centers in 10 Western
120 European countries including Denmark, France, Germany, Greece, Italy, Norway, Spain,
121 Sweden, The Netherlands and the United Kingdom [11]. The design, rationale and methods
122 of the EPIC study including information on dietary assessment methodology, blood collection
123 protocols and follow-up procedures were previously detailed [11].

124 Between 1992 and 1998, standardized lifestyle data, anthropometric measures and biological
125 samples were collected at recruitment, prior to onset of any disease [11]. Validated country-
126 specific questionnaires ensuring high compliance were used to measure diet over the previous
127 12 months [12]. Blood samples are stored at the International Agency for Research on Cancer
128 (IARC, Lyon, France) in -196°C liquid nitrogen for all countries, exceptions being Denmark
129 (nitrogen vapour, -150°C) and Sweden (freezers, -80°C).

130 *The nested case-control study.* The present study focused on data with available sera samples
131 from a nested case-control study in EPIC on hepatocellular carcinoma (HCC) [13]. Cases of
132 HCC were identified from all participating EPIC centres except for Norway and France
133 ($n=117$) from recruitment (1993-1998) up to 2007. Two controls ($n=232$) were selected for
134 each case from all cohort members alive and free of cancer (except non-melanoma skin
135 cancer) by incidence-density sampling and were matched on age at blood collection (± 1 year),
136 sex, study centre, date (± 2 months), time of the day at blood collection (± 3 hours) and fasting
137 status at blood collection (<3 , $3-6$, >6 hours); among women, additional matching criteria
138 included menopausal status (pre-, peri-, post-menopausal) and hormone replacement therapy
139 (HRT) use at time of blood collection (yes/no). In the present study, cases and controls were
140 both included in the analyses as the subjects were all cancer-free at blood collection. Out of
141 the total 349 subjects, 7 subjects (3 cases and 4 controls) had too little serum volume for
142 NMR spectral acquisition with sufficient sensitivity; 6 additional control subjects were
143 excluded following the exclusion of their corresponding case subject. The final analysis
144 included 114 HCC cases and 222 matched controls of which 108 case-control sets with two
145 matched control subjects and 6 sets with one matched control subject.

146 *NMR spectra acquisition.* Sera were processed using standard procedure for ^1H NMR
147 metabolic measurement and profiling protocols [14]. Details on the sera sample preparation as
148 well as NMR data acquisition and processing have been described elsewhere [15]. In brief,

149 each spectrum was reduced to 8,500 bins of 0.001 ppm width over the chemical shift range of
150 0.5 to 9 ppm. Spectra were normalized to total intensity, centred and Pareto scaled, and
151 additionally normalized for batch-effects using the batch profiling calibration method [16].
152 After removal of the structured noise (characterized by a specific mean and standard
153 deviation) located in a well-known noise region (8.5-9ppm) and variables with identical
154 characteristics, the statistical recoupling of variables (SRV) [17], a bucketing procedure, was
155 applied to the metabolomic spectra. The SRV procedure identifies clusters of variables with
156 respect to the ratio of covariance and correlation between consecutive variables along the
157 chemical shift axis, allowing the restauration of the spectral dependency and the recovery of
158 complex NMR signals corresponding to potential physical, chemical or biological entities.
159 More details on the SRV procedure are available in the **Mathematical Appendix**. This
160 permitted a reduction of the number of NMR variables from 8,500 bins to 285 clusters of
161 variables corresponding to reconstructed peak entities which constituted the Y-set of
162 metabolic variables. All steps to obtain the data were done without knowledge of the case-
163 control status of the subjects. Quality control (QC) samples were included to ensure
164 reproducibility of the NMR data acquisition.

165 *Metabolite identification.* The assignment of NMR signals observed in the ^1H one-
166 dimensional fingerprints to metabolites has been achieved by the analysis of additional 2D
167 NMR experiments ^1H - ^{13}C HSQC and ^1H - ^1H TOCSY obtained on a subset of representative
168 samples (one control and one case). The measured chemical shifts were compared to
169 reference shifts of pure compounds using HMDB [18], MMCD [19] and ChenomX,
170 (ChenomX NMR suite, ChenomxInc, Edmonton, Canada) databases.

171 *Lifestyle variables.* The predictors (what will be referred to later on as the X-set) included 13
172 dietary variables from main EPIC food groups compiled from validated country-specific food
173 frequency questionnaires (FFQ) [11,20] (potatoes and other tubers; vegetables; legumes;

174 fruits, nuts and seeds; dairy products; cereal and cereal products; meat and meat products; fish
175 and shellfish; egg and egg products; fat; sugar and confectionary; cakes and biscuits; non-
176 alcoholic beverages), alcohol average lifetime intake (continuous, g/day), anthropometric
177 measures including body mass index (continuous, kg/m²) and height (continuous, cm) that
178 were measured by trained interviewers in the majority of participants [11], highest level of
179 education achieved (categorical: none or primary school completed, technical/professional
180 school, secondary school, longer education (incl. university degree), unspecified), smoking
181 status (categorical: never, former, current smoker, unknown), a measure of physical activity
182 (continuous, metabolic equivalents of task (MET)/h), hepatitis status (yes/no, from biomarker
183 measures of HBV and HCV seropositivity [ARCHITECT HBsAg and anti-HCV
184 chemiluminescent microparticle immunoassays; Abbott Diagnostics, France]) and baseline
185 self-reported diabetes status (yes/no). Descriptive information on these variables can be found
186 in **Supplementary table 1**.

187 *Statistical analysis*

188 *PC-PR2 analysis.* Principal component partial R-square (PC-PR2) was primarily used to
189 identify and quantify sources of systematic variability within metabolomic data [15]. PC-PR2
190 combines aspects of principal component analysis (PCA) and the R^2_{partial} statistic in multiple
191 linear regression, and allows for (some) inter-correlation between the explanatory variables
192 under scrutiny [15]. In short, PCA is performed on the 285 clusters of ¹H NMR variables and
193 a number of components is retained explaining an amount of total variability above a
194 designated threshold (here, 80%). Then, multiple linear regression models are fitted where
195 each component's variability is explained in terms of relevant covariates, e.g. specific
196 characteristics of samples like country of origin, smoking status, laboratory treatment, etc. For
197 each given component, the R^2_{partial} statistic is computed for all covariates, quantifying the
198 amount of variability each independent variable explains, conditional on all other covariates

199 included in the model. Finally, an overall R^2_{partial} is calculated as a weighted average for every
200 covariate, using the eigenvalues as components' weights. Mathematical details pertaining to
201 the PC-PR2 method are described elsewhere [15].

202 In this study, PC-PR2 was applied to the 285 clusters of NMR variables, whereas the
203 explanatory variables examined for systematic variability were NMR batch, country of origin,
204 sex, age at blood collection, serum clot contact time (centrifugation at the day of blood
205 collection d, or the following day, d+1), length of freezing time (≤ 15 vs. >15 years), and
206 fasting status at blood collection (< 3 , 3-6, > 6 hours). With the similar motivation of
207 identifying sources of variability within lifestyle data, a similar PC-PR2 analysis was applied
208 to the 21 lifestyle factors; the examined covariates for systematic variability were country of
209 origin, sex and age at recruitment. For both metabolomics and lifestyle data, residuals on the
210 variable accounting for most variability, identified through PC-PR2 analyses, were computed
211 in a series of univariate linear regression models [21] and were used in the subsequent PLS.

212 *PLS analysis.* A PLS model was used to relate lifestyle variables to metabolomic profiles.

213 PLS is a multivariate technique that generalizes features of PCA and multiple linear
214 regression. PLS iteratively extracts linear combinations of, in turn, predictors (the X-set) and
215 responses (the Y-set), which in this study, were lifestyle variables and metabolomic profiles,
216 respectively. First, components or latent factors are extracted allowing a simultaneous
217 decomposition of the X- and Y-sets, in order to maximize their covariance [22]. The factors
218 extracted from the predictors' set are orthogonal. Computational details of PLS are described
219 in the **Mathematical Appendix**. As a standard step for the PLS algorithm, the X- and Y-sets
220 were centered and standardized for the analysis and a simple expectation-maximization (EM)
221 algorithm, adapted from the PLS kernel algorithm [23,24], was used to compute covariance
222 matrices when missing values were present in the lifestyle data. This was done as follows: a
223 first pass of PLS was computed filling in the missing values by the average of the non-

224 missing values for each corresponding variable. A second pass was then performed whereby
225 the missing data were assigned their predicted values based on the first model, and the PLS
226 regression is recomputed.

227 Then, a seven-fold cross validation analysis was carried out to select the number h of
228 significant PLS factors to retain [8] (see **Mathematical appendix**). This was achieved by
229 splitting the data into seven groups of observations. In turn, each group of observations was
230 considered as the test set, whilst the other six were the training sets, used to perform PLS
231 analysis. A measure of PLS performance was determined for each step through the predicted
232 residual sum of squares (PRESS) statistic, whereby the predicted values in the test set, the \tilde{Y}_h
233 matrix, based on the X-components estimated through the model in the training set, were
234 compared to the observed responses, the Y matrix. This comparison is quantified by the
235 squared Euclidean distance between these two matrices. In turn for an increasing number h of
236 components, the process is iterated seven times, until each group of observations serves as a
237 test set. Eventually, the number h of selected PLS factors is the one minimizing the PRESS
238 statistic.

239 For each PLS factor, loadings were computed for the lifestyle (X-set) and the NMR (Y-set)
240 variables. The loadings, i.e. coefficients quantifying the contribution of each original variable
241 to the PLS factor, were used to characterize the various factors. As the analysis involved
242 many variables in the X-set and, particularly, in the Y-set, the interpretation focused primarily
243 on variables with loading values lower than the 10th percentile and larger than the 90th
244 percentile for the X variables, and lower than the 5th and larger than the 95th percentiles for the
245 Y variables, that were deemed the most significant contributors to the PLS factor.

246 *Logistic regression analysis.* Last, scores of each PLS factor were related to HCC risk in
247 conditional logistic regression models to compute HCC odds ratios (ORs) and associated 95%
248 confidence intervals (95% CI) where ORs express the change in HCC risk associated to one

249 standard deviation (1-SD) increase in the score. Models were adjusted for C-reactive protein
250 concentration, alpha-fetoprotein concentration and for a composite score indicative of liver
251 damage. The score summarizes the number of abnormal values of circulating enzymes
252 measured in the hepatic tissue in six liver function tests (alanine aminotransferase >55 U/L,
253 aspartate aminotransferase >34 U/L, gamma-glutamyltransferase: men>64 U/L and
254 women>36 U/L, alkaline phosphatase >150 U/L, albumin<35 g/L, total bilirubin>20.5
255 $\mu\text{mol/L}$; cut-points were provided by the clinical biochemistry laboratory that conducted the
256 analyses and were based on assay specifications) [25]. These biomarkers were measured on
257 the ARCHITECT c Systems™ and the AEROSSET System (Abbott Diagnostics) using
258 standard protocols. Laboratory analyses were performed at the Centre de Biologie République
259 laboratory, Lyon, France. These adjustments were deemed necessary to address potential
260 confounding stemming from metabolic disorders, inflammation or underlying liver
261 dysfunction [25–28]. Adjustments for total dietary fibre, vitamin D, calcium and iron intakes
262 (continuous) were evaluated but not retained in the final models for lack of confounding
263 exerted by these variables. The receiver operating characteristic (ROC) curve and the
264 associated area under the curve (AUC) were determined from conditional logistic regressions
265 to evaluate the predictive performance of PLS models. AUC values were computed for
266 conditional logistic models including progressively the PLS scores, separately for lifestyle
267 and metabolomic factors (as shown in **Table 4**, column 1). The sensitivity, specificity and
268 accuracy were calculated for a cut-off point, selected as the minimal distance between the
269 ROC curve and the upper left corner of the diagram [29,30]. The corrected positive predictive
270 value (PPV), taking into account the nested case-control design [31,32] was computed by
271 including the prevalence of HCC in the EPIC population ($\pi= 0.0004$), computed over a 7-year
272 period (1992-2010) where 191 HCC cases were ascertained from a total of 477,206
273 participants included for case identification after relevant exclusions [33]. The AUC

274 unavoidably increases with the number of covariates added to the conditional logistic model.
275 To address this issue, a resampling scheme was devised to compute an objective/ unbiased
276 estimate of the AUC, inspired by the work of Uno et al [34]. For each one of the 1000 drawn
277 bootstrap samples, a 10-fold cross-validation was performed, repeated ten times to remove
278 variation due to random partitioning of data and to yield more stable estimates. The predicted
279 values from each of the conditional logistic models in the training set were used to derive
280 AUC values in the test set. The 2.5th and 97.5th percentile values made up the 95% confidence
281 intervals.

282 *Sensitivity analyses.* A sensitivity analysis was performed by running PLS on data excluding
283 sets where cases were diagnosed within the first two years of follow-up. The model was
284 conducted on 271 observations (92 cases, 179 controls), to investigate the performance of the
285 PLS model, ruling out potential reverse causation. The metabolomic profiles of HCC cases
286 diagnosed within two years from enrolment could reflect the presence of the tumour rather
287 than informing about tumour aetiology. The variable importance in the projection (VIP)
288 statistic was used to facilitate the comparison of the sensitivity analysis with the main
289 analysis. The VIP expresses the explanatory power of a predictor variable X across all
290 response variables Y (see **Mathematical Appendix**).

291 *Mediation analysis.* The mediating role of the Y-scores in the association between lifestyle
292 profiles and HCC risk was assessed. Separately for each extracted combination of lifestyle
293 and metabolomic PLS factors, mediation analyses were performed with the ‘paramed’ Stata
294 function that allows for exposure-mediator interaction based on Valeri and VanderWeele’s
295 work [35]. Briefly, mediation was computed using a Baron and Kenny approach adapted to
296 dichotomous outcomes [36], where two models were specified. In the mediator model, the
297 mediator (the Y-score) was linearly regressed on the exposure (the X-score), while in the
298 outcome model the exposure (X-score) and the mediator (Y-score) were related to the HCC

299 indicator in unconditional logistic regressions. Both models accounted for the concentration
300 of C-reactive protein, alpha-fetoprotein and the composite score of liver damage, and
301 additionally accommodated the other extracted metabolic profiles (Y-scores) to control for
302 mediator-outcome confounders that may occur when estimating the Natural Indirect Effect
303 (NIE) [35]. As the outcome (HCC) is rare, direct and indirect effects can be estimated taking
304 into account the case-control design. This is done by using the same formulas for the effects,
305 while running the mediator regression only for the controls [36]. As mediation packages do
306 not yet accommodate conditional logistic models, the outcome and the mediator models,
307 which were accommodated in unconditional logistic regressions, were adjusted for center and
308 age at blood collection for sake of consistency with previous steps of the analysis.

309 Statistical analyses were performed using R [37] and SAS [38] in general, with the following
310 packages for specific purposes: PROC PLS in SAS 9.4 for PLS analyses, ‘paramed’ in Stata
311 12 [39] for mediation analyses, ‘OptimalCutpoints’ in R for ROC-related assessments.

312 The different steps of the analytical framework developed in this study to model the MITM
313 are presented in **Figure 1**.

314 **Results**

315 In the PC-PR2 analyses, a total of 17 and 14 principal components were retained to
316 explain an amount of total variability exceeding 80% in metabolomics and lifestyle data
317 respectively. **Figure 2** shows that the ensemble of explanatory variables accounted for 19.4%
318 and 26.7% of total variance, respectively in metabolomics and lifestyle data, of which the
319 highest contributor was ‘country of origin’ with consistently 8% and 22%. PLS analysis was
320 carried controlling for this variable.

321 After a seven-fold cross-validation, three PLS factors were retained accounting for
322 21.7% and 8.5% of the overall variability observed in predictor and response variables,

323 respectively (**Table 1**). Lifestyle variables and clusters of NMR variables contributing highly
324 to PLS factors were identified using factor loading values (**Table 2**). The first PLS factor was
325 predominantly positively associated with dairy products and cakes and biscuits intake, while
326 lifetime alcohol intake, smoking status and diabetes displayed negative loadings for this
327 lifestyle component (**Table 2**). On the same PLS factor, signals mainly associated with
328 glucose and bonds of lipids with negative loading values, and with aspartate, glutamine and
329 lysine with positive loadings emerged on the metabolomic profile (**Table 2**). Lifestyle
330 variables characterizing the second PLS factor included cereal products, height and education
331 level with negative loadings, and hepatitis with positive loadings. The metabolic signature
332 included NMR variables with positive loadings associated with aromatic amino acids
333 (phenylalanine, tyrosine) and glucose; and those with negative loadings associated mainly
334 with bonds of lipids, threonine and mannose (**Table 2**). The third PLS factor had a lifestyle
335 pattern outlining intake of vegetables (high negative loadings values), lifetime alcohol
336 consumption, smoking, and hepatitis infection (positive loadings). Its counterpart NMR
337 pattern highlighted signals of glucose and aspartate, with high negative loadings, along with
338 signals of ethanol, myo-inositol, proline and glutamate as prominent metabolites with positive
339 loadings (**Table 2**).

340 Conditional logistic regression models relating HCC risk with the X- and Y-scores are
341 shown in **Table 3**. The first PLS factor was associated to a non-significant decreased HCC
342 risk (23% and 4% in the X- and Y-scores respectively), while the second and third factors
343 were associated to a statistically significant increased HCC risk (54% and 11%; and 37% and
344 22% respectively). Results for the ROC curves parameters are reported in **Table 4**, including
345 AUC, sensitivity, specificity, accuracy and PPV for different combinations of the X- and Y-
346 scores. The AUC of the X-scores and Y-scores for all 3 PLS factors, adjusted for C-reactive
347 protein concentration, alpha-fetoprotein concentration and the score of liver damage, was

348 respectively 0.859 and 0.853. An increase in the resampled cross-validated AUC values was
349 also observed for all three X- and Y-scores, albeit smaller, with respectively 0.836 and 0.827.
350 Results from the sensitivity analysis conducted on data excluding sets where cases were
351 diagnosed within the first two years of follow-up, showed similarities in terms of lifestyle
352 variables' and metabolites' loadings on the PLS factors (**Supplementary Table 2**). Notable
353 differences pertained to the identification of new signals for the first PLS factor including
354 ethanol, histidine and an unknown compound. On the second lifestyle factor, BMI (positive
355 loadings) replaced education level (negative loadings) while the reflected metabolomic profile
356 was comparable to its counterpart from the main analysis (**Supplementary Table 2**). On the
357 third factor, smoking status and hepatitis (positive loadings) were replaced by sugar and
358 confectionary intake (negative loadings); signals contributing to the associated metabolic
359 profile remained the same but the direction of the association was inversed as loadings had
360 opposite signs as compared to the counterpart PLS factor of the main model (**Supplementary**
361 **Table 2**). Corresponding ORs from conditional logistic regression models relating the X- and
362 Y-scores to HCC risk are available in **Table 5**. The scores showed a statistically significant
363 association in the second factor for both sets and in the third factor for the Y-set. ROC-
364 associated statistics for different models are presented in **Supplementary Table 3**. The VIP
365 plot (**Figure 3**) displayed the results for the importance of the lifestyle variables in the
366 prediction of the Y-set computed for the main PLS model performed including all subjects
367 (panel **A**) and for the sensitivity model (panel **B**). The results suggested a potential gain in
368 stability as prominent lifestyle variables for prediction were maintained
369 (hepatitis/diabetes/cakes and biscuits), the magnitude of the VIP was improved for some
370 (fat/lifetime alcohol intake) and less emphasis was put on others (BMI/physical activity).
371 Finally, the natural indirect effect was assessed in the mediation analyses and the results are
372 presented in **Table 6**. Overall, there was limited evidence that metabolomic signals mediated

373 the association between lifestyle components and HCC risk in the first PLS factor. Evidence
374 of a significant mediated effect by the Y-scores was found in the second and third PLS factors
375 when models were adjusted for exposure-mediator interaction (**Table 6**).

376 **Discussion**

377 In this work, an analytical strategy based on PLS analysis was conceived to extract
378 relevant information from sets of lifestyle and NMR metabolomic variables, and to relate the
379 resulting components to the risk of disease. This offered a way to implement the MITM
380 approach [6] in a nested case-control study on HCC within the EPIC study. MITM has been
381 suggested as a way to link specific putative metabolites to lifestyle exposures and disease
382 outcomes, thus leading to the identification of potential intermediate biomarkers [6].

383 An implementation of MITM was previously carried out in a nested case-control study
384 in the Turin sub-cohort of EPIC [7] based on prospectively collected plasma samples from a
385 pilot study on colon and breast cancers. In their work, a list of intermediate markers was
386 identified by an in-parallel evaluation of the relationships between untargeted ¹H NMR
387 profiles with dietary exposures and risk of colon and breast cancers using correlation analysis
388 and logistic regression. In our study, a different analytical framework was developed, largely
389 exploiting features of PLS analysis, a multivariate technique that iteratively extracts
390 components capturing co-variability in sets of predictors and response variables [8,40]. A set
391 of lifestyle predictor variables were related to NMR responses. In a second step, PLS
392 predictors' and responses' scores were linked to the risk of HCC.

393 Another sensitive issue in this analysis was the choice of lifestyle variables. Two
394 disease-indicator variables reflecting environmental exposures, diabetes and hepatitis, were
395 included in the set of predictors, as they turned out to have an important role in the
396 characterization of metabolomic signatures. In addition, diabetes is the main metabolic risk

397 factor for HCC alongside with fatty liver disease [41,42], and chronic infection with hepatitis
398 B (HBV) and particularly hepatitis C (HCV) viruses were classified as class I carcinogens for
399 HCC by IARC [43].

400 Other relevant biomarkers were not part of the list of predictors in PLS analysis, but were
401 controlled for in logistic regression models. This included C-reactive protein, alpha-
402 fetoprotein, and a score for liver damage, an index of different circulating enzymes measured
403 in the hepatic tissue indicating potential underlying liver function impairment [25]. The alpha-
404 fetoprotein was included as an adjustment factor in the analyses not because of its established
405 part as a serum marker for HCC diagnosis [26,44], but rather to account for it as a potential
406 confounder that may cloud the relation between scores and HCC, both in conditional logistic
407 regressions and in mediation analyses.

408 Similarly to other multivariate techniques, a key aspect of PLS analysis is the choice
409 of the number of factors to retain, in an effort of exhaustively summarizing data variability
410 through a limited number of factors. Based on a seven-fold cross-validation, three linear
411 combinations of variables were extracted in this work. A challenging aspect of this analysis is
412 the interpretation of these factors, with respect to lifestyle and metabolomic variables. A
413 subjective criterion based on the distribution of loading values was used throughout. The
414 variables displaying the most extreme loading values (in absolute terms) were the ones
415 characterizing each factor.

416 The first lifestyle factor highlighted a healthy pattern with negative loadings for
417 diabetes status, smoking status and lifetime alcohol intake, and was not associated to HCC
418 risk, similarly to its metabolomics counterpart. The lifestyle component of the second PLS
419 factor, was reflective of a lifestyle pattern reflective of “higher-risk exposures”, and was
420 related to a significant 54% increase in HCC risk. Likewise, its associated metabolic
421 component displayed a significant HCC risk augmentation by 11%. The lifestyle component

422 of the third PLS factor described participants with lower vegetables intake, elevated lifetime
423 alcohol consumption, more likely to be ever smokers and hepatitis positive; one standard
424 deviation increase of this component was associated to a statistically significant 37% increase
425 in HCC risk. Similarly, a 22% significant increase in HCC risk was observed for its metabolic
426 counterpart, characterized by positive signals of ethanol and myo-inositol, and displayed
427 negative loadings for glucose.

428 The MITM is captured by the rationale of PLS analysis, in the sense that each set of lifestyle
429 profiles and metabolic signatures of the extracted PLS factors mirrored one another. In
430 addition, mediation was observed for the second and third PLS factors, whereby the
431 metabolomic component mediated the relation between the lifestyle component and HCC, for
432 which statistically significant associations with HCC risk were estimated, emphasising the
433 presence of a MITM. Mediation analysis relies on the assumption that there is no mediator-
434 outcome confounder that is affected by the exposure [35]. In our study C-reactive protein,
435 alpha-fetoprotein and liver damage score were weakly correlated to lifestyle factor score, thus
436 introducing potential bias in the estimation of direct and indirect effects in our mediation
437 analysis. Additionally, a number of background confounders (mediator-outcome and
438 exposure-outcome confounders) were present that we have tried to control for, either by
439 adjustments or by accounting for potential interactions, however some degree of bias can
440 remain and caution should be employed when interpreting the results.

441 The predictive performance of PLS factors in relation to HCC occurrence was evaluated
442 through an analysis of AUC values. The performance of the model improved progressively,
443 with all 3 X- and Y-scores added; after a bootstrapped cross-validation, the AUC estimates
444 were lower but the increase in the performance was nevertheless present. The ROC
445 methodology allows estimation of PPV, which expresses the risk of disease after a positive
446 test [45]. In a setting with low HCC prevalence ($\pi=0.0004$), in line with Western populations

447 [46], extremely low PPV estimates were observed. In the absence of a very specific test, many
448 positive tests arise from disease-free individuals [45], thus leading to a dilution of PPV.
449 A sensitivity analysis was carried out excluding the first two years of follow-up, but results
450 were virtually unchanged, both in terms of relative risk estimates in logistic regression
451 models, and of percentage of variability explained in PLS analysis. These findings suggest
452 that reverse causation bias, if present, was minimal.

453 This study had the ambition of integrating in the same analytical framework study
454 participants' lifestyle characteristics with a large number of NMR metabolic profiles. These
455 data pose a number of methodological challenges due to their size and the complexity of
456 exhaustively capturing and interpreting the biological processes they reflect. To address these
457 issues, techniques involving multivariate statistics have been progressively revived in the
458 recent years [2]. Epidemiologic evaluations of metabolomic data frequently combined PLS
459 with discriminant analysis, such as PLS-DA or O-PLS-DA. The main objective of these
460 methods is to identify a series of metabolomic features distinguishing between two very
461 distinct groups of study participants [47,48]. In such strategies, only one set of variables is
462 multi-dimensional and the response is one variable only. Similar multivariate techniques for
463 pattern extraction, belonging to the family of regression methods, include reduced rank
464 regression. This multivariate method relates an ensemble of response variables to a set of
465 predictor variables where the estimated matrix of the regression coefficients is of reduced
466 rank [49–51]. In addition, canonical correlation analysis (CCA) [52] is a method applied to
467 identify the optimum structure or dimensionality of each variable set that maximizes the
468 relationship between two sets of multi-dimensional variables. The main difference between
469 CCA and PLS regression is that CCA maximizes the correlation between the two new
470 dimensions, i.e. extracted factors, whereas PLS maximizes their covariance. PLS can be
471 considered as a trade-off between CCA and PCA, since maximizing the covariance

472 corresponds to maximizing the product of the correlation and standard deviation, given that
473 $\text{cov}(X,Y)=\text{cor}(X,Y)*\text{SD}(X)*\text{SD}(Y)$.

474 Untargeted NMR was used in this work to acquire metabolomic signals. Prior to PLS
475 analysis, a bucketing procedure, the statistical recoupling of variables (SRV) [17,53], was
476 applied to reduce the number of NMR variables to 285 clusters. This was done by aggregating
477 consecutive NMR bins based on their covariance to correlation ratio. This allowed the
478 identification of informative components of the spectra, thus acting as an efficient noise-
479 removing filter. Subsequently the annotation effort remains challenging, for a number of
480 reasons. The majority of published metabolomics studies often identified a limited number of
481 metabolites at a time [54], and the Human Metabolome Database (HMDB) and other related
482 resources [18,55], that offer richly annotated information continuously increasing the
483 metabolite coverage for users, are mostly exploited through time consuming interactive
484 procedures. In addition, individual metabolites often overlap in NMR signals, which can
485 hinder annotations. These challenges, as well as large variability in metabolite concentrations,
486 and disentangling informative signals from noise, are not specific to NMR and pertain to any
487 type of untargeted technique. Such investigations may profit from complementary targeted
488 metabolomic analytical strategies [55].

489 Throughout the different steps of this work, the scaling problem was first tackled by
490 normalizing spectra to total intensity. NMR data were also centered and Pareto-scaled,
491 together with correction for potential batch effects [16]. The PC-PR2 method offered a way to
492 investigate major sources of systematic variability in NMR and lifestyle data [15]. The
493 variable “country of origin” emerged as the variable accounting for the largest proportion of
494 total variability, and the residual method was used to control for this variable in the following
495 steps of the analysis. While this may lead to removing regional gradients of dietary
496 variability, this step is instrumental to avoid unwanted systematic regional-specific bias in the

497 data in country-specific questionnaire assessments. In addition, technical aspects like storage
498 and handling of biological samples, fasting status at blood collection are specific to each
499 country [15]. In any case, variability due to “country of origin” is not exploited in conditional
500 logistic models, as cases and controls were also matched on center.

501 One of the limitations of this study is the restricted sample size which raises concerns
502 with regards to power to detect associations. While a larger sample size would possibly result
503 in more statistically significant findings, we used the data that was available with NMR
504 profiles measured. In this work we have developed a framework to analyse complex data
505 integrating lifestyle and metabolomics in relation to risk of disease. The approach described in
506 this study has merits but also pitfalls among which it is worth mentioning that statistical
507 methods are used repeatedly on the same set of data, notably the PLS model, the conditional
508 logistic regression, the AUC estimation and mediation analysis. To partially address this, a
509 cross-validation approach was devised for AUC estimation which involved conditional
510 logistic regression, whereby PLS was done without knowledge of the case/control status.
511 However, conditional logistic regression models and mediation analyses were implemented
512 on the same data, and our analysis did not account for this limitation. This may have led to
513 spuriously increase the nominal level of statistical significance of statistical tests.

514 **Conclusion**

515 The MITM emerged as a method for the identification of relevant biomarkers, with
516 great potential to unravel utmost important steps in the aetiology of disease. The analytical
517 strategy for MITM was developed to use all potentially informative aspects of high-
518 throughput data by integrating metabolomic, dietary and lifestyle exposures together with
519 disease indicators. While the framework was applied towards the investigation of HCC
520 determinants, it can be easily extended to similar aetiological contexts and applied to other –
521 omics settings.

522 **Funding**

523 This work was supported by the French National Cancer Institute (L'Institut National du
524 Cancer; INCA) [grant number 2009-139; PI: M. Jenab]. The coordination of the European
525 Prospective Investigation into Cancer and nutrition is financially supported by the European
526 Commission (Directorate General for Health and Consumer Affairs) and the International
527 Agency for Research on Cancer. The national cohorts are supported by Health Research Fund
528 (FIS) of the Spanish Ministry of Health RTICC 'Red Temática de Investigación Cooperativa
529 en Cáncer [Grant numbers: Rd06/0020/0091, Rd12/0036/0018], Regional Governments of
530 Andalucía, Asturias, Basque Country, Murcia [project 6236] and Navarra, Instituto de Salud
531 Carlos III, Redes de Investigación Cooperativa (RD06/0020) (Spain); the Danish Cancer
532 Society (Denmark); the Ligue Contre le Cancer, the Institut Gustave Roussy, the Mutuelle
533 Générale de l'Éducation Nationale, and the Institut National de la Santé et de la Recherche
534 Médicale (France); the Deutsche Krebshilfe, the Deutsches Krebsforschungszentrum, and the
535 Federal Ministry of Education and Research (Germany); the Hellenic Health Foundation, the
536 Stavros Niarchos Foundation, and the Hellenic Ministry of Health and Social Solidarity
537 (Greece); the Italian Association for Research on Cancer AIRC and the National Research
538 Council (Italy); the Dutch Ministry of Public Health, Welfare and Sports, the Netherlands
539 Cancer Registry, LK Research Funds, Dutch Prevention Funds, the Dutch Zorg Onderzoek
540 Nederland, the World Cancer Research Fund, and Statistics Netherlands (Netherlands);
541 European Research Council-2009-AdG 232997 and the Nordforsk, Nordic Centre of
542 Excellence program on Food, Nutrition and Health (Norway); the Swedish Cancer Society,
543 the Swedish Research Council, and the Regional Governments of Skåne and Västerbotten
544 (Sweden); Cancer Research UK, the Medical Research Council, the Stroke Association, the
545 British Heart Foundation, the Department of Health, the Food Standards Agency, and the
546 Wellcome Trust (United Kingdom). The work undertaken by N Assi was supported by by the
547 Université de Lyon I through a doctoral fellowship awarded by the EDISS doctoral school.

548 **Acknowledgments**

549 We would like to acknowledge the assistance of Dr Elodie Jobard from the ISA-CRMN in
550 obtaining the annotation of the NMR data.

551 **Reference List**

- 552 1. Nicholson, J.K., Holmes, E., and Elliott, P. (2008) The metabolome-wide association
553 study: a new look at human disease risk factors. *J. Proteome Res.*, **7**,3637–3638.
- 554 2. Chadeau-Hyam, M., Campanella, G., Jombart, T., Bottolo, L., Portengen, L., Vineis,
555 P., Liquet, B., and Vermeulen, R.C.H. (2013) Deciphering the Complex:
556 Methodological Overview of Statistical Models to Derive OMICS-Based Biomarkers.
557 *Environ. Mol. Mutagen.*, **54**,542–557.
- 558 3. Floegel, A., Wientzek, A., Bachlechner, U., Jacobs, S., Drohan, D., Prehn, C.,
559 Adamski, J., Krumsiek, J., Schulze, M.B., Pischon, T., and Boeing, H. (2014) Linking
560 diet, physical activity, cardiorespiratory fitness and obesity to serum metabolite
561 networks: findings from a population-based study. *Int. J. Obes. (Lond)*, **1**–9.
- 562 4. Trushina, E. and Mielke, M.M. (2013) Recent advances in the application of
563 metabolomics to Alzheimer’s Disease. *Biochim. Biophys. Acta*, **1842**,1232–1239.
- 564 5. Jin, X., Yun, S.J., Jeong, P., Kim, I.Y., Kim, W.-J., and Park, S. (2014) Diagnosis of
565 bladder cancer and prediction of survival by urinary metabolomics. *Oncotarget*,
566 **5**,1635–1645.
- 567 6. Vineis, P. and Perera, F. (2007) Molecular epidemiology and biomarkers in etiologic
568 cancer research: the new in light of the old. *Cancer Epidemiol. Biomarkers Prev.*,
569 **16**,1954–1965.
- 570 7. Chadeau-Hyam, M., Athersuch, T.J., Keun, H.C., De Iorio, M., Ebbels, T.M.D., Jenab,
571 M., Sacerdote, C., Bruce, S.J., Holmes, E., and Vineis, P. (2011) Meeting-in-the-
572 middle using metabolic profiling - a strategy for the identification of intermediate
573 biomarkers in cohort studies. *Biomarkers*, **16**,83–88.
- 574 8. Tenenhaus, M. (1998) La régression PLS. Technip. Paris.
- 575 9. Mitra, V. and Metcalf, J. (2009) Metabolic functions of the liver. *Anaesth. Intensive*
576 *Care Med.*, **10**,334–335.
- 577 10. Fages, A., Duarte-Salles, T., Stepien, M., Ferrari, P., Fedirko, V., Pontoizeau, C.,
578 Trichopoulou, A., Aleksandrova, K., Tjønneland, A., Olsen, A., Clavel-Chapelon, F.,
579 Boutron-Ruault, M.-C., Severi, G., Kaaks, R., Kuhn, T., Floegel, A., Boeing, H.,
580 Lagiou, P., Bamia, C., Trichopoulos, D., Palli, D., Pala, V., Panico, S., Tumino, R.,
581 Vineis, P., Bueno-de-Mesquita, H.B., Peeters, P.H.M., Weiderpass, E., Agudo, A.,
582 Molina-Montes, E., Huerta, J.M., Ardanaz, E., Dorronsoro, M., Sjöberg, K., Ohlsson,
583 B., Khaw, K.-T., Wareham, N.J., Travis, R.C., Schmidt, J.A., Cross, A.J., Gunter, M.J.,
584 Riboli, E., Scalbert, A., Romieu, I., Elena-Herrmann, B., and Jenab, M. (2015)
585 Metabolomic Profiles of Hepatocellular Carcinoma in a European Prospective Cohort.
586 *Submitt. to Am. J. Gastroenterol.*,.
- 587 11. Riboli, E., Hunt, K.J., Slimani, N., Ferrari, P., Norat, T., Fahey, M., Charrondière,
588 U.R., Hémon, B., Casagrande, C., Vignat, J., Overvad, K., Tjønneland, A., Clavel-
589 Chapelon, F., Thiébaud, A., Wahrendorf, J., Boeing, H., Trichopoulos, D.,

- 590 Trichopoulou, A., Vineis, P., Palli, D., Bueno-De-Mesquita, H.B., Peeters, P.H.M.,
591 Lund, E., Engeset, D., González, C. a, Barricarte, A., Berglund, G., Hallmans, G., Day,
592 N.E., Key, T.J., Kaaks, R., and Saracci, R. (2002) European Prospective Investigation
593 into Cancer and Nutrition (EPIC): study populations and data collection. *Public Health*
594 *Nutr.*, **5**,1113–1124.
- 595 12. Kaaks, R., Slimani, N., and Riboli, E. (1997) Pilot Phase Studies on the Accuracy of
596 Dietary Intake Measurements in the EPIC Project : Overall Evaluation of Results.
597 **26**,26–36.
- 598 13. Trichopoulos, D., Bamia, C., Lagiou, P., Fedirko, V., Trepo, E., Jenab, M., Pischon, T.,
599 Nöthlings, U., Overved, K., Tjønneland, A., Outzen, M., Clavel-Chapelon, F., Kaaks,
600 R., Lukanova, A., Boeing, H., Aleksandrova, K., Benetou, V., Zylis, D., Palli, D., Pala,
601 V., Panico, S., Tumino, R., Sacerdote, C., Bueno-De-Mesquita, H.B., Van Kranen,
602 H.J., Peeters, P.H.M., Lund, E., Quirós, J.R., González, C. a, Sanchez Perez, M.-J.,
603 Navarro, C., Dorronsoro, M., Barricarte, A., Lindkvist, B., Regnér, S., Werner, M.,
604 Hallmans, G., Khaw, K.-T., Wareham, N., Key, T., Romieu, I., Chuang, S.-C.,
605 Murphy, N., Boffetta, P., Trichopoulou, A., and Riboli, E. (2011) Hepatocellular
606 carcinoma risk factors and disease burden in a European cohort: a nested case-control
607 study. *J. Natl. Cancer Inst.*, **103**,1686–1695.
- 608 14. Beckonert, O., Keun, H.C., Ebbels, T.M.D., Bundy, J., Holmes, E., Lindon, J.C., and
609 Nicholson, J.K. (2007) Metabolic profiling, metabolomic and metabonomic procedures
610 for NMR spectroscopy of urine, plasma, serum and tissue extracts. *Nat. Protoc.*,
611 **2**,2692–2703.
- 612 15. Fages, A., Ferrari, P., Monni, S., Dossus, L., Floegel, A., Mode, N., and Al., E. (2014)
613 Investigating sources of variability in metabolomic data in the EPIC study: the
614 Principal Component Partial R-square (PC-PR2) method. *Metabolomics*, **10**,1074–
615 1083.
- 616 16. Fages, A., Pontoizeau, C., Jobard, E., Lévy, P., Bartosch, B., and Elena-Herrmann, B.
617 (2013) Batch profiling calibration for robust NMR metabonomic data analysis. *Anal.*
618 *Bioanal. Chem.*, **405**,8819–8827.
- 619 17. Blaise, B.J., Shintu, L., Elena, B., Emsley, L., Dumas, M.-E., and Toulhoat, P. (2009)
620 Statistical Recoupling Prior to Significance Testing in Nuclear Resonance Based
621 Metabonomics. *Anal. Chem.*, **81**,6242–6251.
- 622 18. Wishart, D.S., Knox, C., Guo, A.C., Eisner, R., Young, N., Gautam, B., Hau, D.D.,
623 Psychogios, N., Dong, E., Bouatra, S., Mandal, R., Sinelnikov, I., Xia, J., Jia, L., Cruz,
624 J. a, Lim, E., Sobsey, C. a, Shrivastava, S., Huang, P., Liu, P., Fang, L., Peng, J.,
625 Fradette, R., Cheng, D., Tzur, D., Clements, M., Lewis, A., De Souza, A., Zuniga, A.,
626 Dawe, M., Xiong, Y., Clive, D., Greiner, R., Nazyrova, A., Shaykhutdinov, R., Li, L.,
627 Vogel, H.J., and Forsythe, I. (2009) HMDB: a knowledgebase for the human
628 metabolome. *Nucleic Acids Res.*, **37**,D603–D610.
- 629 19. Cui, Q., Lewis, I.A., Hegeman, A.D., Anderson, M.E., Li, J., Schulte, C.F., Westler,
630 W.M., Eghbalnia, H.R., Sussman, M.R., and Markley, J.L. (2008) Metabolite

- 631 identification via the Madison Metabolomics Consortium Database. *Nat. Biotechnol.*,
632 **26**,162–164.
- 633 20. Slimani, N., Deharveng, G., Unwin, I., Southgate, D.A., Viignat, J., Skeie, G., Salvini,
634 S., Parpinel, M., Moller, A., and et al. (2007) The EPIC nutrient database project
635 (ENDB): a first attempt to standardize nutrient databases across the 10 European
636 countries participating in the EPIC study - DTU Orbit. *Eur. J. Clin. Nutr.*, **61**,1037–
637 1056.
- 638 21. Kleinbaum, D.G., Kupper, L.K., and Muller, K.E. (1987) Applied regression analysis
639 and other multivariable methods. Belmont, CA: Duxbury Press.
- 640 22. Wold, S., Sjostrom, M., and Ericksson, L. (2001) PLS-regression : a basic tool of
641 chemometrics. *Chemom. Intell. Lab. Syst.*, **58**,109–130.
- 642 23. Rannar, S., Geladi, P., Lindgren, F., and Wold, S. (1995) A PLS kernel algorithm for
643 data sets with many variables and few objects. Part II: cross-validation, missing data
644 and examples. *J. Chemom.*, **9**,459–470.
- 645 24. Bastien, P. (2008) Régression PLS et Données Censurées. Conservatoire National des
646 Arts et Métiers - CNAM.
- 647 25. Fedirko, V., Trichopolou, A., Bamia, C., Duarte-Salles, T., Trepo, E., Aleksandrova,
648 K., Nöthlings, U., Lukanova, A., Lagiou, P., Boffetta, P., Trichopoulos, D., Katzke, V.
649 a, Overvad, K., Tjønneland, A., Hansen, L., Boutron-Ruault, M.C., Fagherazzi, G.,
650 Bastide, N., Panico, S., Grioni, S., Vineis, P., Palli, D., Tumino, R., Bueno-de-
651 Mesquita, H.B., Peeters, P.H., Skeie, G., Engeset, D., Parr, C.L., Jakszyn, P., Sánchez,
652 M.J., Barricarte, A., Amiano, P., Chirlaque, M., Quirós, J.R., Sund, M., Werner, M.,
653 Sonestedt, E., Ericson, U., Key, T.J., Khaw, K.T., Ferrari, P., Romieu, I., Riboli, E.,
654 and Jenab, M. (2013) Consumption of fish and meats and risk of hepatocellular
655 carcinoma: the European Prospective Investigation into Cancer and Nutrition (EPIC).
656 *Ann. Oncol.*, **24**,2166–2173.
- 657 26. Akuta, N., Suzuki, F., Kobayashi, M., Hara, T., Sezaki, H., Suzuki, Y., Hosaka, T.,
658 Kobayashi, M., Saitoh, S., Ikeda, K., and Kumada, H. (2014) Correlation Between
659 Hepatitis B Virus Surface antigen Level and Alpha-Fetoprotein in Patients Free of
660 Hepatocellular Carcinoma or Severe Hepatitis. *J. Med. Virol.*, **86**,131–138.
- 661 27. Kanazir, M., Boricic, I., Delic, D., Tepavcevic, D.K., Knezevic, A., Jovanovic, T., and
662 Pekmezovic, T. (2010) Risk factors for hepatocellular carcinoma: A case-control study
663 in Belgrade (Serbia). *Tumori*, **96**,911–917.
- 664 28. Zheng, Z., Zhou, L., Gao, S., Yang, Z., Yao, J., and Zheng, S. (2013) Prognostic role of
665 C-reactive protein in hepatocellular carcinoma: A systematic review and meta-analysis.
666 *Int. J. Med. Sci.*, **10**,653–664.
- 667 29. Metz, C.D. (1978) Basic Principles of ROC Analysis. *Semin. Nucl. Med.*, **8**,283–298.

- 668 30. Vermont, J., Bosson, J.L., François, P., Robert, C., Rueff, A., and Demongeot, J.
669 (1991) Strategies for graphical threshold determination. *Comput. Methods Programs*
670 *Biomed.*, **35**,141–150.
- 671 31. Biesheuvel, C.J., Vergouwe, Y., Oudega, R., Hoes, A.W., Grobbee, D.E., and Moons,
672 K.G.M. (2008) Advantages of the nested case-control design in diagnostic research.
673 *BMC Med. Res. Methodol.*, **8**,48.
- 674 32. Van Zaane, B., Vergouwe, Y., Donders, a R.T., and Moons, K.G.M. (2012)
675 Comparison of approaches to estimate confidence intervals of post-test probabilities of
676 diagnostic test results in a nested case-control study. *BMC Med. Res. Methodol.*,
677 **12**,166.
- 678 33. Stepien, M., Duarte-Salles, T., Fedirko, V., Floegel, A., Kumar-Barupal, D., Rinaldi,
679 S., Achaintre, D., Assi, N., Tjønneland, A., Overvad, K., Bastide, N., Boutron-Ruault,
680 M.-C., Severi, G., Kuhn, T., Kaaks, R., Aleksandrova, K., Boeing, H., Trichopoulou,
681 A., Bamia, C., Lagiou, P., Saieva, C., Agnoli, C., Panico, S., Tumino, R., Naccarati, A.,
682 Bueno-de-Mesquita, H.B., Peeters, P.H., Weiderpass, E., Quirós, J.R., Agudo, A.,
683 Sanchez, M.-J., Dorronsoro, M., Gavrila, D., Barricarte, A., Ohlsson, B., Sjöberg, K.,
684 Werner, M., Sund, M., Wareham, N., Khaw, K.-T., Travis, R.C., Schmidt, J.A., Gunter,
685 M., Cross, A.J., Vineis, P., Romieu, I., Scalbert, A., and Jenab, M. (2015) Alteration of
686 Amino Acid and Biogenic Amine Metabolism in Hepatobiliary Cancers: Findings from
687 a Prospective Cohort Study. *Submitt. to Int. J. Cancer.*,
- 688 34. Uno, H., Cai, T., Pencina, M.J., D'Agostino, R.B., and Wei, L.J. (2011) On the C-
689 statistics for Evaluating Overall Adequacy of Risk Prediction Procedures with
690 Censored Survival Data. *Stat. Med.*, **30**,1105–1117.
- 691 35. Valeri, L. and Vanderweele, T.J. (2013) Mediation analysis allowing for exposure-
692 mediator interactions and causal interpretation: theoretical assumptions and
693 implementation with SAS and SPSS macros. *Psychol. Methods*, **18**,137–150.
- 694 36. Vanderweele, T.J. and Vansteelandt, S. (2010) Odds ratios for mediation analysis for a
695 dichotomous outcome. *Am. J. Epidemiol.*, **172**,1339–1348.
- 696 37. R Foundation for Statistical Computing and R Core Team. (2013) R: A language and
697 environment for statistical computing.
- 698 38. SAS Institute Inc. and Cary, N. (2012) Base SAS® 9.4 Procedures Guide.
- 699 39. College Station TX : StataCorp. (2011) Stata Statistical Software: Release 12.
- 700 40. Abdi, H. (2010) Partial least squares regression and projection on latent structure
701 regression (PLS Regression). *Wiley Interdiscip. Rev. Comput. Stat.*, **2**,97–106.
- 702 41. Yang, W.-S., Va, P., Bray, F., Gao, S., Gao, J., Li, H.-L., and Xiang, Y.-B. (2011) The
703 role of pre-existing diabetes mellitus on hepatocellular carcinoma occurrence and
704 prognosis: a meta-analysis of prospective cohort studies. *PLoS One*, **6**,e27326.

- 705 42. Gomaa, A.-I. (2008) Hepatocellular carcinoma: Epidemiology, risk factors and
706 pathogenesis. *World J. Gastroenterol.*, **14**,4300–4308.
- 707 43. Cogliano, V.J., Baan, R., Straif, K., Grosse, Y., Lauby-Secretan, B., El Ghissassi, F.,
708 Bouvard, V., Benbrahim-Tallaa, L., Guha, N., Freeman, C., Galichet, L., and Wild,
709 C.P. (2011) Preventable exposures associated with human cancers. *J. Natl. Cancer*
710 *Inst.*, **103**,1827–1839.
- 711 44. Bialecki, E.S. and Di Bisceglie, A.M. (2005) Diagnosis of hepatocellular carcinoma.
712 *HPB (Oxford)*., **7**,26–34.
- 713 45. Wentzensen, N. and Wacholder, S. (2013) From differences in means between cases
714 and controls to risk stratification: a business plan for biomarker development. *Cancer*
715 *Discov.*, **3**,148–157.
- 716 46. Leong, T.Y.-M. and Leong, A.S.-Y. (2005) Epidemiology and carcinogenesis of
717 hepatocellular carcinoma. *HPB (Oxford)*., **7**,5–15.
- 718 47. Rothwell, J., Fillâtre, Y., Martin, J.-F., Lyan, B., Pujos-Guillot, E., Fezeu, L., Hercberg,
719 S., Comte, B., Galan, P., Touvier, M., and Manach, C. (2014) New biomarkers of
720 coffee consumption identified by the non-targeted metabolomic profiling of cohort
721 study subjects. *PLoS One*, **9**,e93474.
- 722 48. Guo, M., Zhao, B., Liu, H., Zhang, L., Peng, L., Qin, L., Zhang, Z., Li, J., Cai, C., and
723 Gao, X. (2014) A Metabolomic Strategy to Screen the Prototype Components and
724 Metabolites of Shuang-Huang-Lian Injection in Human Serum by Ultra Performance
725 Liquid Chromatography Coupled with Quadrupole Time-of-Flight Mass Spectrometry.
726 *J. Anal. Methods Chem.*, **2014**,241505.
- 727 49. Anderson, T.W. (1951) Estimating linear restrictions on regression coefficients for
728 multivariate normal distributions. *Ann. Math. Stat.*, **22**,327–351.
- 729 50. Izenman, A.J. (1975) Reduced-Rank Regression for the Multivariate Linear Model. *J.*
730 *Multivar. Anal.*, **5**,248–264.
- 731 51. Aldrin, M. (2002) Reduced-rank regression. In: El-Shaarawi AH, Piegorisch WW,
732 editors. Encyclopedia of Environmetrics. Chichester: John Wiley & Sons, Ltd. pp.
733 1724–1728.
- 734 52. Hotelling, H. (1936) Relations between two sets of variates. *Biometrika*, **28**,321–377.
- 735 53. Navratil, V., Pontoizeau, C., Billoir, E., and Blaise, B.J. (2013) SRV : an opensource
736 toolbox to accelerate the recovery of metabolic biomarkers and correlations from
737 metabolic phenotyping data sets. *Bioinformatics*, **29**,1348–1349.
- 738 54. Wishart, D.S. (2008) Quantitative metabolomics using NMR. *TrAC Trends Anal.*
739 *Chem.*, **27**,228–237.
- 740 55. Psychogios, N., Hau, D.D., Peng, J., Guo, A.C., Mandal, R., Bouatra, S., Sinelnikov, I.,
741 Krishnamurthy, R., Eisner, R., Gautam, B., Young, N., Xia, J., Knox, C., Dong, E.,

742 Huang, P., Hollander, Z., Pedersen, T.L., Smith, S.R., Bamforth, F., Greiner, R.,
743 McManus, B., Newman, J.W., Goodfriend, T., and Wishart, D.S. (2011) The human
744 serum metabolome. *PLoS One*, **6**,e16957.

745

746 **Legends to figures**

747 **Figure 1:** General scheme of the analytical framework developed in the study. A PC-PR2
748 analysis is carried out beforehand to identify relevant sources of variation. In the PLS model
749 the X- and Y- sets are related to each other, and scores are computed (1). X- and Y-scores are,
750 in turn, associated to a case-control indicator of HCC status in conditional logistic regression
751 models (2). A mediation analysis is carried out to explore the role of metabolomics in the
752 association between lifestyle factors and risk of HCC (3).

753 **Figure 2:** PC-PR2 analysis results* identifying the sources of variability in the NMR data
754 (panel A) and in the lifestyle data (panel B).

755 * 17 and 14 components were retained to account for 80 % (threshold used) of total NMR (A)
756 and lifestyle variability (B), respectively. The R2 value represents the amount of variability in
757 NMR / lifestyle variable explained by the ensemble of investigated predictors.

758 **Figure 3:** Variable importance plot (VIP) displaying the variable importance for projection
759 statistic of the predictor variables for the PLS analyses.

760 Panel A: Results from the main PLS model run on all observations (N=336, X-set=21, Y-
761 set=285).

762 Panel B: Results from the PLS sensitivity analysis run on a subsample (N=271, 92 cases, 179
763 controls) excluding sets where cases were diagnosed within the first two years of follow-up
764 (X-set=21, Y-set=285).

765 The horizontal line corresponds to Wold's criterion (0.8), the threshold used to rule if a
766 variable has an important contribution to the construction of the Y variables (see
767 **Mathematical Appendix** for further details).

Table 1: Individual and cumulative variation (%) explained by the first 3 PLS factors in 21 lifestyle (X-set) and 285 NMR (Y-set) variables.

# of PLS Factors	Lifestyle Variables		NMR Variables	
	Individual	Cumulative	Individual	Cumulative
1	6.17	-	5.51	-
2	6.23	12.40	2.38	7.89
3	9.27	21.67	0.59	8.48

Table 2: Lifestyle and NMR cluster variables contributing to each of the 3 PLS factors (N=336, X-set=21, Y-set=285).

PLS Factor	Lifestyle Variable*	Loading value	CS*‡ (ppm)	Metabolite**	Loading value
1	Dairy Products	0.28	5.22		-0.06
	Cakes and Biscuits	0.32	3.88		-0.05
	Lifetime Alcohol Consumption	-0.25	3.82		-0.06
	Smoking Status	-0.39	3.76		-0.06
	Diabetes	-0.63	3.71	Glucose	-0.05
			3.54		-0.05
			3.50		-0.07
			3.48		-0.07
			3.44	Acetoacetate	-0.07
			3.23	Choline + Glycerphosphocholine	-0.04
			3.01	Lysine	0.10
			2.94	Albumin	0.10
			2.65	Aspartate	0.10
			2.42	Glutamine	0.10
			2.28	Acetoacetate	0.10
			2.22	CH ₂ -CH ₂ -COOC bond of lipids + Acetone	-0.04
			1.86		0.09
			1.87	Lysine	0.10
			1.53	CH ₂ -CH ₂ -COOC bond of lipids	-0.03
	2	Cereal and Cereal Products	-0.16	7.17	Tyrosine
Height		-0.34	6.87		0.13
Education Level		-0.26	5.27	CH=CH bond of lipids	-0.13
Hepatitis		0.49	5.22	Glucose	0.16
			5.18	Mannose + Lipid O-CH ₂	-0.12
			4.27	Lipid O-CH ₂	-0.12
			4.25	Threonine	-0.14
			4.07	Choline + Lipid O-CH ₂ + Myo-inositol	-0.12
			4.05	Creatinine	-0.14
			3.88		0.15
			3.82		0.16
			3.76		0.15
			3.71	Glucose	0.15
			3.54		0.15
			3.50		0.16
			3.48		0.16
			3.44	Acetoacetate	0.16
			3.23	Choline + Glycerphosphocholine	0.15
			2.80	Aspartate	-0.12
			2.22	CH ₂ -CH ₂ -COOC bond of lipids + Acetone	-0.11
		2.19	CH ₂ -CH ₂ -COOC bond of lipids	-0.15	
		2.02	Proline + Glutamate + CH ₂ =C bonds of lipids	-0.13	
		1.53	CH ₂ -CH ₂ -COOC bond of lipids	-0.13	
		1.25	CH ₂ bond of lipids	-0.12	
		0.86	Cholesterol + CH ₃ bond of lipids	-0.12	
3	Vegetables	-0.42	7.32	Phenylalanine	0.11
	Lifetime Alcohol Consumption	0.29	5.22	Glucose	-0.13
	Smoking Status	0.25	4.28	Lipid O-CH ₂	0.11
	Hepatitis	0.26	3.88		-0.11
			3.82		-0.11
			3.76	Glucose	-0.12
			3.71		-0.11
			3.69		-0.11
			3.63	Myo-inositol	0.16
			3.50		-0.13
			3.48	Glucose	-0.12
			3.44	Acetoacetate	-0.12
			3.35		0.11
			3.33	Proline	0.13
			3.28	Myo-inositol	0.12
			3.23	Choline + Glycerphosphocholine	-0.12
			2.80	Aspartate	-0.13
			2.76	part of =CH-CH ₂ -CH= bond of lipids	-0.13
			2.35		0.12
			2.33	Proline + Glutamate	0.13

	1.20	3-hydroxybutyrate + CH2 bond of lipids	0.11
	1.16	Ethanol	0.15
	0.66	Cholesterol	0.11

*Relevant lifestyle and NMR variables contributing to each PLS factor selected based on their associated loading values <10th percentile (pctl) and >90th pctl or <5th pctl and >95th pctl respectively.

‡ CS: ¹H chemical shift (in ppm) of the cluster (center value).

**Some of the identified clusters were found to be background noise during the annotation phase and were removed from this table.

Table 3: HCC odds ratios* and 95% confidence interval (OR, 95% CI) associated with the lifestyle (X-set) and the NMR clusters (Y-set) PLS scores in the main analysis (N=336, X-set=21, Y-set=285).

PLS Lifestyle Variables X-scores			PLS NMR Variables Y-scores		
Factor	OR** (95% CI)	P-Wald†	Factor	OR** (95% CI)	P-Wald†
1	0.77 (0.58, 1.02)	0.07	1	0.96 (0.91, 1.01)	0.09
2	1.54 (1.06, 2.25)	0.02	2	1.11 (1.02, 1.22)	0.02
3	1.37 (1.05, 1.79)	0.02	3	1.22 (1.04, 1.44)	0.01

*Models were adjusted for C-reactive protein concentration, alpha-fetoprotein concentration and a composite score for liver damage. Cases and controls were matched on age at blood collection (± 1 year), sex, study centre, date (± 2 months) and time of the day at blood collection (± 3 hours), fasting status at blood collection (<3/3-6/>6 hours); among women, additional matching criteria included menopausal status (pre-/peri-/postmenopausal) and hormone replacement therapy use at time of blood collection (yes/no). ** ORs expressing the change in HCC risk associated to 1-SD increase in the score. † Wald's test was for continuous exposure compared with a Chi-square distribution with 1 degree of freedom (dof).

Table 4: Area under the curve (AUC), sensitivity, specificity, accuracy and positive predictive value (PPV) of ROC models (with 95% CI), from the main PLS analysis (N=336, X-set=21, Y-set=285).

	AUC	AUC _b **	Sensitivity	Specificity	Accuracy	PPV
Adjustment Covariates (ADJ)*	0.842 (0.794, 0.891)	0.821 (0.766, 0.868)	0.752 (0.662, 0.829)	0.802 (0.743, 0.852)	0.785	0.0015
X1 scores + ADJ	0.846 (0.797, 0.894)	0.825 (0.766, 0.875)	0.743 (0.653, 0.821)	0.838 (0.783, 0.884)	0.806	0.0018
X1+X2 scores + ADJ	0.854 (0.808, 0.900)	0.831 (0.772, 0.881)	0.743 (0.653, 0.821)	0.824 (0.768, 0.872)	0.797	0.0017
X1+X2+X3 scores + ADJ	0.859 (0.811, 0.907)	0.836 (0.778, 0.887)	0.796 (0.710, 0.866)	0.788 (0.729, 0.840)	0.791	0.0015
Y1 scores + ADJ	0.841 (0.793, 0.890)	0.817 (0.760, 0.865)	0.735 (0.643, 0.813)	0.820 (0.763, 0.868)	0.791	0.0016
Y1+Y2 scores + ADJ	0.845 (0.795, 0.894)	0.820 (0.762, 0.872)	0.735 (0.643, 0.813)	0.851 (0.798, 0.895)	0.812	0.0020
Y1+Y2+Y3 scores + ADJ	0.853 (0.804, 0.902)	0.827 (0.771, 0.877)	0.726 (0.634, 0.805)	0.883 (0.833, 0.922)	0.890	0.0025

*The model is run on the adjustment covariates (ADJ) including the C-reactive protein concentration, alpha-fetoprotein concentration and a composite score for liver damage. ** AUC_b is the bootstrapped-cross validated estimate of the AUC. X1, X2 and X3 are the lifestyle component scores of the first, second and third PLS factors, respectively. Y1, Y2, and Y3 are the metabolomics component of the first, second and third PLS factors, respectively.

Table 5: HCC odds ratios* and 95% confidence intervals (OR, 95%CI) associated with the lifestyle (X-set) and the NMR clusters (Y-set) PLS scores. Results from the sensitivity analysis (N=271, 92 cases, 179 controls) conducted excluding sets where cases were diagnosed within the first two years of follow-up (X-set=21, Y-set=285).

PLS Lifestyle Variables X-scores			PLS NMR Variables Y-scores		
Factor	OR** (95% CI)	P-Wald†	Factor	OR** (95% CI)	P-Wald†
1	0.80 (0.60, 1.08)	0.15	1	0.96 (0.94, 1.04)	0.56
2	1.56 (1.02, 2.40)	0.04	2	1.18 (1.03, 1.36)	0.02
3	0.86 (0.67, 1.11)	0.26	3	0.86 (0.73, 0.99)	<0.05

*Models were adjusted for C-reactive protein concentration, alpha-fetoprotein concentration and a composite score for liver damage. Cases and controls were matched on age at blood collection (± 1 year), sex, study centre, date (± 2 months) and time of the day at blood collection (± 3 hours), fasting status at blood collection (<3/3-6/>6 hours); among women, additional matching criteria included menopausal status (pre-/peri-/postmenopausal) and hormone replacement therapy use at time of blood collection (yes/no). ** ORs expressing the change in HCC risk associated to 1-SD increase in the score. † Wald's test was for continuous exposure compared with a Chi-square distribution with 1 degree of freedom (dof).

Table 6: Results from the mediation analysis (N= 336, X-set=21, Y-set=285): Natural Indirect Effect (NIE) and 95%CI*.

Model**				Natural Indirect Effect (NIE)	
Exposure (A)	Mediator (M)	Outcome	A*M interaction term	Estimate (95%CI)	p-value
X1 score	Y1 score	HCC	No	0.91 (0.77, 1.06)	0.23
X2 score	Y2 score	HCC	No	1.11 (0.97, 1.25)	0.12
X3 score	Y3 score	HCC	No	1.08 (0.94, 1.23)	0.28
X1 score	Y1 score	HCC	Yes	0.96 (0.79, 1.17)	0.70
X2 score	Y2 score	HCC	Yes	1.15 (1.01, 1.31)	0.04
X3 score	Y3 score	HCC	Yes	1.13 (1.01, 1.28)	0.04

* The standard errors used to compute the 95%CI were obtained using the delta method.

**Models were adjusted for the C-reactive protein concentration, alpha-fetoprotein concentration and a composite score for liver damage, as well as for the other Y-scores, as potential mediator-outcome confounders. Additionally, the outcome and the mediator models were adjusted for centre and age at blood collection.

Supplementary Tables

A Statistical framework to model the meeting-in-the-middle principle using metabolomic data: application to hepatocellular carcinoma in the EPIC study.

Supplementary Table 1: Summary statistics of the predictors variables (X-set) of the study subjects in the EPIC liver nested case–control study (N=336, 114 Cases, 222 Controls).

	Mean / N*	sd / %*	p5	p95	N missing
Dietary Variables (g/day)					
Potatoes and other tubers	100.57	78.15	9.34	266.97	0
Vegetables	194.20	143.22	45.03	473.45	0
Legumes	9.85	18.03	0.00	41.18	0
Fruits, nuts and seeds	232.80	197.94	23.55	585.22	0
Dairy products	334.40	261.46	49.92	777.48	0
Cereal and cereal products	227.04	117.67	76.39	458.94	0
Meat and meat products	115.97	62.29	37.83	236.32	0
Fish and shellfish	32.88	32.26	3.78	81.43	0
Egg and egg products	18.67	18.72	1.88	55.57	0
Fat	34.61	18.48	11.01	70.76	0
Sugar and confectionary	47.26	51.51	1.93	138.73	0
Cakes and biscuits	41.33	49.68	0.00	147.26	0
Non-alcoholic beverages	1053.91	793.31	85.00	2391.90	0
Anthropometric variables					
BMI (kg/m ²)	27.41	4.41	21.22	36.16	0
Height (cm)	169.70	9.99	152.00	184.80	0
Lifestyle Variables					
Lifetime alcohol intake (g/day)	23.27	41.38	0	91.998	61
Physical activity (Mets/h)	77.13	49.45	11.5	173.63	0
Highest Education Level					
None or primary school completed	167	49.7	-	-	-
Technical/professional school	75	22.32	-	-	-
Secondary school	27	8.04	-	-	-
Longer education (incl. university degree)	62	18.45	-	-	-
Unspecified or Unknown	5	1.49	-	-	-
Smoking status					
Never	124	36.9	-	-	-
Former	125	37.2	-	-	-
Current smoker	85	25.3	-	-	-
Unspecified or Unknown	2	0.6	-	-	-
Pathology variables indicative of lifestyle					
Hepatitis status					
No	291	86.87	-	-	-
Yes	44	13.13	-	-	-
Diabetes					
No	307	91.37	-	-	-
Yes	29	8.63	-	-	-

*Mean and standard deviation (sd), were reported for continuous variables and frequencies and percentages (%) were reported for categorical variables.

p5: 5th percentile, p95:95th percentile.

Supplementary Table 2: Results from the sensitivity analysis run on a subsample (N=271, 92 cases, 179 controls) excluding sets where cases were diagnosed within the first two years of follow-up (X-set=21, Y-set=285). Lifestyle and NMR cluster variables contributing to each PLS factor.

PLS Factor	Lifestyle Variable*	Loading value	CS*‡ (ppm)	Metabolite**	Loading value
1	Dairy Products	0.33	7.03	Histidine	0.09
	Cakes and Biscuits	0.34	5.22		-0.07
	Lifetime Alcohol Consumption	-0.34	3.88		-0.06
	Smoking Status	-0.26	3.82		-0.07
	Diabetes	-0.59	3.76	Glucose	-0.06
			3.71		-0.06
			3.54		-0.05
			3.50		-0.07
			3.48		-0.08
			3.44	Acetoacetate	-0.08
			3.23	Choline + Glycerphosphocholine	-0.05
			3.03	Creatine	0.10
			3.01	Albumin	0.10
			2.28	Acetoacetate	0.10
			2.22	CH ₂ -CH ₂ -COOC bond of lipids + Acetone	-0.03
			2.06	Proline + Glutamate	0.09
			1.91	Lysine + Arginine	-0.03
			1.87	Lysine	0.09
			1.16	Ethanol	-0.04
			1.08	Unknown 1	0.09
		0.91	CH ₃ bond of lipids	0.09	
2	Cereal and Cereal Products	-0.24	7.17	Tyrosine	0.14
	BMI	0.34	6.87		0.14
	Height	-0.39	5.27	CH=CH bond of lipids	-0.14
	Hepatitis	0.55	5.22	Glucose	0.13
			5.18	Mannose + Lipid O-CH ₂	-0.13
			4.27	Lipid O-CH ₂	-0.12
			4.25	Threonine	-0.14
			4.05	Creatinine	-0.14
			3.88		0.13
			3.82		0.13
			3.76		0.13
			3.75	Glucose	0.12
			3.71		0.12
			3.54		0.15
			3.50		0.13
			3.48		0.13
			3.44	Acetoacetate	0.13
			3.23	Choline + Glycerphosphocholine	0.12
			2.80	Aspartate	-0.13
			2.76	=CH-CH ₂ -CH= bond of lipids	-0.12
		2.19	CH ₂ -CH ₂ -COOC bond of lipids	-0.16	
		2.02	Proline + Glutamate	-0.14	
		1.53	CH ₂ -CH ₂ -COOC bond of lipids	-0.13	
		1.25	CH ₂ bond of lipids	-0.12	
		0.86	Cholesterol + CH ₃ bond of lipids	-0.12	
3	Vegetables	0.39	5.25	Glucose	0.17
	Sugar and Confectionnary	-0.21	4.28	Lipid O-CH ₂	-0.07
	Lifetime Alcohol Consumption	-0.29	4.14	Proline	-0.08
			4.07	Choline + Lipid O-CH ₂ + Myo-inositol	-0.07
			3.88		0.16
			3.82		0.16
			3.76	Glucose	0.16
			3.75		0.14
			3.71		0.15
			3.69		0.16
		3.63	Myo-inositol	-0.16	

	3.54		0.12
	3.50	Glucose	0.17
	3.48		0.17
	3.44	Acetoacetate	0.16
	3.41		-0.10
	3.35	Proline	-0.15
	3.34		-0.12
	3.28	Myo-inositol	-0.09
	3.23	Choline + Glycerphosphocholine	0.15
	1.91	Lysine + Arginine	-0.07
	1.16	Ethanol	-0.16
	0.68		-0.06
	0.66	Cholesterol	-0.08

*Relevant lifestyle and NMR variables contributing to each PLS factor selected based on their associated loading values <10th percentile (pctl) and >90th pctl or <5th pctl and >95th pctl respectively.

‡ CS: ¹H chemical shift (in ppm) of the cluster (center value).

**Some of the identified clusters were found to be background noise during the annotation phase and were removed from this table.

Supplementary Table 3: Results from the sensitivity analysis (N=271, 92 cases, 179 controls) conducted excluding sets where cases were diagnosed within the first two years of follow-up (X-set=21, Y-set=285). Area under the curve (AUC), sensitivity, specificity, accuracy and positive predictive value (PPV) of ROC models (with 95% CI).

	AUC	AUC _b **	Sensitivity	Specificity	Accuracy	PPV
Adjustment Covariate (ADJ) [†]	0.846 (0.793, 0.899)	0.827 (0.765, 0.879)	0.750 (0.649, 0.834)	0.838 (0.776, 0.889)	0.808	0.0018
X1 scores + ADJ	0.853 (0.800, 0.905)	0.834 (0.774, 0.890)	0.728 (0.626, 0.816)	0.872 (0.813, 0.917)	0.823	0.0023
X1+X2 scores + ADJ	0.860 (0.811, 0.910)	0.837 (0.772, 0.893)	0.750 (0.649, 0.834)	0.832 (0.769, 0.884)	0.804	0.0018
X1+X2+X3 scores + ADJ	0.861 (0.810, 0.912)	0.837 (0.773, 0.893)	0.761 (0.661, 0.844)	0.838 (0.776, 0.889)	0.812	0.0019
Y1 scores + ADJ	0.847 (0.794, 0.900)	0.827 (0.768, 0.884)	0.739 (0.637, 0.825)	0.838 (0.776, 0.889)	0.804	0.0018
Y1+Y2 scores + ADJ	0.848 (0.794, 0.901)	0.827 (0.764, 0.883)	0.717 (0.614, 0.806)	0.899 (0.846, 0.939)	0.838	0.0028
Y1+Y2+Y3 scores + ADJ	0.853 (0.800, 0.907)	0.826 (0.763, 0.882)	0.717 (0.614, 0.806)	0.911 (0.859, 0.948)	0.845	0.0032

*The model is run on the adjustment covariates (ADJ) including the C-reactive protein concentration, alpha-fetoprotein concentration and a composite score for liver damage. ** AUC_b is the bootstrapped-cross validated estimate of the AUC. X1, X2 and X3 are the lifestyle component scores of the first, second and third PLS factors, respectively. Y1, Y2, and Y3 are the metabolomics component of the first, second and third PLS factors, respectively.