

A family of partial-linear single-index models for analyzing complex environmental exposures with continuous, categorical, time-to-event, and longitudinal health outcomes

Yuyan Wang

NYU Langone Health

Yinxiang Wu

NYU Langone Health

Melanie H. Jacobson

NYU Langone Health

Myeonggyun Lee

NYU Langone Health

Peng Jin

NYU Langone Health

Leonardo Trasande

NYU Langone Health

Mengling Liu (✉ mengling.liu@nyulangone.org)

NYU Langone Health <https://orcid.org/0000-0001-9758-8522>

Research

Keywords: Environmental mixtures, NHANES, Semiparametric model, Triglyceride

Posted Date: July 14th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-20209/v2>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published on September 11th, 2020. See the published version at <https://doi.org/10.1186/s12940-020-00644-4>.

1 **Title**

2 A family of partial-linear single-index models for analyzing complex environmental exposures
3 with continuous, categorical, time-to-event, and longitudinal health outcomes
4

5 **Author names and affiliations**

6 **Yuyan Wang¹, Yinxiang Wu¹, Melanie H. Jacobson³, Myeonggyun Lee¹, Peng Jin¹,**
7 **Leonardo Trasande^{1,2,3}, Mengling Liu^{1,2,*}**

8 ¹Department of Population Health, NYU Langone Health, New York, NY, USA

9 ²Department of Environmental Medicine, NYU Langone Health, New York, NY, USA

10 ³Department of Pediatrics, Divisions of Environmental Pediatrics, NYU Langone Health, New
11 York, NY, USA

12 Yuyan Wang: yuyan.wang@nyulangone.org;

13 Yinxiang Wu: yinxiang.wu@nyulangone.org;

14 Melanie Jacobson: melanie.jacobson2@nyulangone.org;

15 Myeonggyun Lee: myeonggyun.lee@nyulangone.org;

16 Peng Jin: peng.jin@nyulangone.org;

17 Leonardo Trasande: leonardo.trasande@nyulangone.org;

18

19 ***Corresponding author**

20 **Mengling Liu, PhD**

21 Email address: mengling.liu@nyulangone.org

22 Postal address: 180 Madison Avenue, New York, NY, 10016

23 Tel: 646-501-3652

24 **Abstract**

25 **Background:** Statistical methods to study the joint effects of environmental factors are of great
26 importance to understand the impact of correlated exposures that may act synergistically or
27 antagonistically on health outcomes. This study proposes a family of statistical models under a
28 unified partial-linear single-index (PLSI) modeling framework, to assess the joint effects of
29 environmental factors for continuous, categorical, time-to-event, and longitudinal outcomes. All
30 PLSI models consist of a linear combination of exposures into a single index for practical
31 interpretability of relative direction and importance, and a nonparametric link function for
32 modeling flexibility.

33 **Methods:** We presented PLSI linear regression and PLSI quantile regression for continuous
34 outcome, PLSI generalized linear regression for categorical outcome, PLSI proportional hazards
35 model for time-to-event outcome, and PLSI mixed-effects model for longitudinal outcome.
36 These models were demonstrated using a dataset of 800 subjects from NHANES 2003-2004
37 survey including 8 environmental factors. Serum triglyceride concentration was analyzed as a
38 continuous outcome and then dichotomized as a binary outcome. Simulations were conducted to
39 demonstrate the PLSI proportional hazards model and PLSI mixed-effects model. The
40 performance of PLSI models was compared with their counterpart parametric models.

41 **Results:** PLSI linear, quantile, and logistic regressions showed similar results that the 8
42 environmental factors had both positive and negative associations with triglycerides, with a-
43 Tocopherol having the most positive and trans-b-carotene the most negative association. For the
44 time-to-event and longitudinal settings, simulations showed that PLSI models could correctly
45 identify directions and relative importance for the 8 environmental factors. Compared with
46 parametric models, PLSI models got similar results when the link function was close to linear,
47 but clearly outperformed in simulations with nonlinear effects.

48 **Conclusions:** We presented a unified family of PLSI models to assess the joint effects of
49 exposures on four commonly-used types of outcomes in environmental research, and
50 demonstrated their modeling flexibility and effectiveness, especially for studying environmental
51 factors with mixed directional effects and/or nonlinear effects. Our study has expanded the
52 analytical toolbox for investigating the complex effects of environmental factors. A practical
53 contribution also included a coherent algorithm for all proposed PLSI models with R codes
54 available.

55 **Keywords:** Environmental mixtures, NHANES, Semiparametric model, Triglyceride

56 **Background**

57 Humans are constantly exposed to a mixture of environmental factors that have the potential to
58 affect health adversely or beneficially, such as chemical contaminants, air pollutants, dietary
59 factors, and behavioral and socioeconomic characteristics. The *exposome*, which is defined as the
60 totality of environmental (non-genetic) exposures from conception onwards (i.e., environmental
61 factors), has been proposed to address the complexities related to studying multiple exposures
62 (1). It is well acknowledged that single-exposure-outcome approaches do not allow for the
63 disentangling of effects of multiple exposures, and miss the interplay among them (2). Therefore,
64 quantifying the complex effects of multiple and simultaneous environmental exposures on health
65 outcomes has become a focus of environmental health research (3, 4). The National Institute of
66 Environmental Health Sciences (NIEHS) has been supporting and conducting combined
67 exposure research, and highlighted this direction as a priority in its 2018-2023 Strategic Plan (5).

68 Statistical approaches have been proposed to assess the effects of multiple exposures on
69 health outcomes from different perspectives, each focusing on distinct scientific questions (2, 6).
70 However, several challenges for statistical modeling are apparent in these investigations (2).

71 First, multiple environmental exposures occur simultaneously, often with complex correlation
72 structures among them. Second, they may exhibit synergistic or antagonistic effects on the health
73 outcome, and their associations with health outcomes can be positive, negative, or null, which
74 reflect the complex web of physiological relationships and/or “reverse causality” (7, 8). Third,
75 the relationships between environmental factors and health outcomes can be non-linear, which
76 pose challenges to standard parametric regression-based methods (9). Fourth, it is well
77 recognized that statistical methods have different strengths in addressing various aspects of
78 scientific investigations. For example, from the methodology perspective, Stafoggia et al (2)
79 classified the statistical methods for analysis of environmental mixtures into dimension
80 reduction, variable selection, or grouping or clustering. From the view of scientific questions,
81 Gibson et al (4) distinguished different study objectives as: identifying the important components
82 in the mixtures, studying synergistic effects, or characterizing the overall effect of the mixtures.

83 Specifically, in studying the joint effects of environmental exposures, weighted quantile sum
84 regression (WQS) (9, 10) and Bayesian kernel machine regression (BKMR) (11, 12) are two
85 popular modeling approaches. The WQS method is a parametric method assuming that all
86 exposures are associated with the outcome in one direction in each run of analysis, and then
87 derives a one-dimensional weighted sum score of the exposures under the assumed direction for
88 the estimation of overall effect. BKMR is a nonparametric method and can handle nonlinear and
89 complex relationships between exposure mixtures and outcome. Some measures have been
90 proposed to quantify the importance and effects of exposure components based on BKMR
91 results. For example, the posterior inclusion probability (PIP) characterizes the probability of an
92 exposure being associated with outcome, and change per interquartile range increase quantifies
93 the expected change in the outcome in association with the change in an exposure from the 25th
94 to 75th percentile, while other exposures are fixed to the median. However, the nonparametric

95 exposure-response function may be difficult to interpret and its fitting often needs a large sample
96 size (13, 14). In addition, WQS and BKRM have been generalized to study environmental
97 mixtures with several types of outcomes, such as WQS for longitudinal outcomes (15) and
98 BKMR for time-to-event outcomes (16). However, a general modeling framework that can
99 alleviate the above limitations in environmental health research is still desired (17).

100 Partial-linear single-index (PLSI) models are a family of semiparametric models that reside
101 between the completely unstructured nonparametric models and restrictive parametric regression
102 models (18-20). By reducing multiple exposures into a single index through a linear combination
103 of the exposures, the PLSI models can reduce the “curse of dimensionality” issue and improve
104 modeling efficiency. The application and performance of single-index linear regression for
105 analysis of environmental exposures with continuous outcomes has been evaluated previously
106 (pending publication). Specifically, the PLSI modeling framework allows the associations
107 between exposures and outcomes to be in the positive or negative direction, provides explicit and
108 interpretable quantification on the relative direction and importance of the exposures, and models
109 these effects with flexibility through a nonparametric link function. Therefore, PLSI models are
110 able to address the objectives of identifying important individual exposures, their direction and
111 magnitude of association with the outcome, and characterizing the overall effect of multiple
112 exposures or exposure mixtures, responding well to the key scientific objectives summarized by
113 Gibson et al. (4). In recent years, research on PLSI models has attracted increasing attention and
114 extended to different types of outcomes, such as categorical (21-23), time-to-event (24-27) and
115 longitudinal (28-31) outcomes. Table 1 summarizes the outcome types of interest and
116 corresponding PLSI models with key references and their corresponding counterpart parametric
117 models.

118 **Table 1** Summary of outcome types and corresponding PLSI models and parametric models

Outcome type	PLSI models	Counterpart models	Key references	Equation
Continuous	PLSI linear regression	Linear regression	(18), (21), (22), (32),(33), (34), (35), (36), (37), (38)	(1)
	PLSI quantile regression	Quantile regression	(39), (40), (41), (42), (43), (44)	(2)
Categorical (binary)	PLSI generalized linear (logistic) regression	Generalized linear (logistic) regression	(18), (22), (36), (38)	(3)
Time-to-event	PLSI PH model	Cox PH model	(24), (25), (26), (27)	(4)
Longitudinal	PLSI mixed-effects model	Linear mixed-effects model	(28), (29), (45), (46), (47)	(5)

119 The main goal of this study was to unify the resource advantages of PLSI models into one
120 general framework for analyzing environmental factors, and to demonstrate their values in
121 environmental research for different types of health outcomes. We exemplified the use of PLSI
122 models in assessing the associations between correlated environmental factors with health
123 outcomes using real and simulated datasets based on National Health and Nutrition Examination
124 Survey (NHANES) 2003-2004 cycle. Another aim was to develop effective computation
125 algorithms for the PLSI models and to consolidate these models using R packages.

126 **Methods**

127 **NHANES dataset**

128 To demonstrate the PLSI models, we used the data from the NHANES 2003-2004 cycle based on
129 the original paper by Patel et al (48), which systematically evaluated the associations of
130 environmental factors with serum lipid levels. We used serum triglyceride concentrations as the
131 primary outcome for demonstration and also considered three demographic variables, age, sex,
132 and race/ethnicity as potential confounders. Participants with data on serum triglycerides,
133 environmental factors and confounders were included in this study (n=800). Details on data pre-
134 processing are provided in [Additional file 1: Figure S1](#). Subjects provided written informed
135 consent, and the Institutional Review Board of the National Center for Health Statistics approved
136 the survey (49). [Table 2](#) summarizes the final variables included in analyses, and [Figure 1](#) shows

137 the correlation matrix of the final 8 environmental factors and triglycerides. The dataset is
 138 provided as [Additional file 2](#), and the R codes conducting data cleaning is included in the R
 139 markdown file ([Additional file 3](#)).

140 **Table 2** List of analyzed variables from NHANES 2002-2003 dataset

Type	Variable name	Abbreviations	Symbol
Outcome	Triglycerides (mg/dL)	TG	Y
Environmental factors	a-Tocopherol (ug/dL)	a-Tocopherol	X1
	g-tocopherol (ug/dL)	g-tocopherol	X2
	Retinyl palmitate (ug/dL)	Retinyl-palmitate	X3
	Retinol (ug/dL)	Retinol	X4
	3,3',4,4',5-Pentachlorobiphenyl (pncb) Lipid Adj (pg/g)	3,3,4,4,5-pncb	X5
	Polychlorinated Biphenyl (PCB) 194 Lipid Adj (ng/g)	PCB156	X6
	2,3,4,6,7,8-hxcdf Lipid Adj (pg/g)	2,3,4,6,7,8-hxcdf	X7
	trans-b-carotene (ug/dL)	trans-b-carotene	X8
Confounders	Age (years)	Age	Z1
	Sex (1: male; 2: female)	Sex	Z2
	Race/Ethnicity (1: Non-Hispanic white; 2: Non-Hispanic black; 3: Mexican American; 4: Other race - Including multi-racial; 5: Other Hispanic)	Race	Z3

141 **Notation and PLSI models overview**

142 For notational convention throughout this article, we let Y denote the outcome, $X = (X_1, \dots, X_8)$
 143 denote the 8 exposure variables to be modeled into the “single index” term, and vector Z
 144 represent the confounders (age, sex, and race/ethnicity). The outcome, continuous triglycerides,
 145 and all exposure variables, except for retinol, were log-transformed, and all exposure variables
 146 were standardized to have mean of zero and standard deviations of 1 before model fitting.

147 In contrast to standard generalized linear models (GLMs) that specify the effects of
 148 exposures and confounders all linearly as $\beta'X + \gamma'Z$, PLSI models assume the influence of
 149 exposures X through a nonparametric link function on the single index while modeling other
 150 confounders linearly, i.e. $g(\beta'X) + \gamma'Z$. The single index coefficients β 's characterize the

151 relative direction and importance of each exposure X_i , and γ for the corresponding linear
 152 coefficient vector for confounder vector Z . Because the link function $g(\cdot)$ is completely
 153 nonparametric, to ensure model identifiability, the l_2 norm of β 's (i.e. $\sqrt{\beta_1^2 + \dots + \beta_8^2}$) is set to be
 154 1 with the first component $\beta_1 > 0$, which are the commonly used parametrization constraints for
 155 all PLSI models (22, 36-38). PLSI models are not identifiable without these constraints because
 156 any scaling or constant shift can be absorbed by the nonparametric link function.

157 **Continuous outcome: mean regression**

158 The PLSI linear regression model is considered as a generalization of both standard linear
 159 regression and missing-link function problem in linear modeling (50), and specified as

$$160 \quad Y = g\left(\sum_{j=1}^8 \beta_j X_j\right) + \gamma'Z + \varepsilon \quad (1)$$

161 The semiparametric PLSI linear regression has the parametric component $\sum_{j=1}^8 \beta_j X_j$ and $\gamma'Z$ for
 162 easy linear representation and interpretation, and the nonparametric components $g(\cdot)$ is totally
 163 unspecified and represents the overall effect of single index, which incorporates potential
 164 nonlinearity and interactions among exposures. When the estimated $g(\cdot)$ is monotone, the effect
 165 of X_j can be interpreted qualitatively using the sign of β_j . If $g(\cdot)$ is monotone increasing, then a
 166 positive sign for β_j suggests increased conditional expectation of Y at larger value of X_j , and vice
 167 versa for a negative sign. As the overall scale of β is set, $|\beta_j|$ can be explained as the relative
 168 importance of X_j affecting the mean of outcome Y as X_j is perturbed while $g(\cdot)$ and other
 169 variables are held fixed. We can also intuitively interpret β_j^2 as the proportion of contribution to
 170 the single index by variable X_j because, when (X_1, X_2, \dots, X_8) are independent, β_j^2 simply
 171 represents X_j 's variance contribution.

172 Besides the analysis for the 8 selected exposures, we also conducted a sensitivity analysis
173 including all 22 environmental factors to investigate the performance of PLSI linear regression to
174 handle highly correlated exposures (Additional file 1: Figure S2).

175 Continuous outcome: quantile regression

176 Beyond the commonly-considered effects of environmental factors on the mean of a continuous
177 outcome, sometimes we are interested in the specific relations cross multiple points of the
178 outcome's distribution, such as higher quantiles of triglycerides (51), higher quantiles of blood
179 pressure (52), low quantiles of birth weight (53), or lower quantiles of intelligence quotient
180 scores (54). Moreover, when the distribution of continuous outcome deviates from Gaussian,
181 modeling the median can be more robust than evaluating the mean by conventional linear
182 regression (55). For this purpose, quantile regression (QR), which was originally proposed by
183 Koenker and Bassett (56) and used as a useful technique in econometrics (57) and growth curve
184 analysis (58), enables us to study the associations of environmental factors with continuous
185 health outcomes as various quantiles across its distribution. PLSI quantile regression is a
186 combination of the PLSI technique and QR (42, 43), and thus we consider it for the analysis of
187 joint effects of multiple environmental factors on the quantile(s) of continuous outcome variable.

188 Given a specific $\tau \in (0,1)$, the PLSI quantile regression for the τ th conditional quantile θ_τ of
189 continuous outcome Y given environmental factors X and covariates Z can be specified as

$$190 \quad \theta_\tau(Y|X, Z) = g_\tau \left(\sum_{j=1}^8 \beta_{\tau j} X_j \right) + \gamma'_\tau Z \quad (2)$$

191 Interpretation of coefficients β_τ 's in the PLSI quantile regression is similar to that of PLSI linear
192 regression, with the difference being that the associations are now with the conditional quantiles
193 of outcome variable $\theta_\tau(Y|X, Z)$ instead of the mean.

194 Categorical outcome: generalized linear regression

195 PLSI generalized linear regression can be employed for categorical outcomes, such as binary,
196 multinomial, or count variables. Here we considered the binary outcome of high triglycerides (>
197 150 milligrams per deciliter) (59), which accounted for 30.75% of the 800 subjects. The PLSI
198 logistic model is specified as

$$199 \quad \text{logit}(P(Y = 1|X, Z)) = g\left(\sum_{j=1}^8 \beta_j X_j\right) + \gamma'Z \quad (3)$$

200 The interpretation of coefficients is based on the log odds that response value is '1' conditioning
201 on the predictors, and β_j represents the relative direction and importance of X_j associated with
202 the log odds of high triglycerides when scale of β is set and $g(\cdot)$ and other variables are held
203 fixed. The logit function can be adapted accordingly to the type of categorical outcome, and the
204 model specifications for multinomial and count outcomes were provided in Additional file 1:
205 Table S1.

206 Time-to-event outcome: proportional hazards model

207 The Cox proportional hazards (PH) regression has been the pivotal model in time-to-event
208 analysis since Sir Cox proposed it in 1972 (60, 61). The Cox PH regression models the hazard
209 function and assumes that covariates have linear effects on the log hazard function. Combining
210 PLSI modeling technique and Cox PH regression, the PLSI PH model is specified as

$$211 \quad \lambda(t|X, Z) = \lambda_0(t) \exp\left\{g\left(\sum_{j=1}^8 \beta_j X_j\right) + \gamma'Z\right\}, \quad (4)$$

212 where β_j can be explained as the relative effect direction and importance of X_j on the log hazard
213 function and $g(\cdot)$ characterizes the overall effect of the index.

214 Longitudinal outcome: mixed-effects model

215 Longitudinal studies arise frequently in environmental research, in which outcomes are measured
216 repeatedly over a period of time with either baseline or time-dependent environmental factors.
217 As measurements from the same subject are often correlated, subject-specific random effects are
218 used to accommodate within-subject dependence and to explain across-subject heterogeneity.
219 Mixed-effects models provide a general and flexible framework for modeling longitudinal data,
220 consisting of two modeling components: fixed effects and random effects, characterizing the
221 population mean and individual variation, respectively (62, 63). Mixed-effects models in general
222 are amenable to missing data and can accommodate missing completely at random or missing at
223 random (62, 64). Without loss of generality, we consider a longitudinal study with N subjects
224 and the i th subject has n_i observations over time. Repeated measures of the outcome are denoted
225 by Y_{ij} , exposure vector X_{ij} , covariate vector Z_{ij} and observation time T_{ij} , and then the observed
226 full dataset is $\{(Y_{ij}, X_{ij}, Z_{ij}, T_{ij}), i = 1, \dots, N, j = 1, \dots, n_i\}$.

227 Specifically, the PLSI mixed-effects model with a random intercept is specified as

$$228 \quad Y_{ij} = g\left(\sum_{l=1}^8 \beta_l X_{ijl}\right) + Z'_{ij}\gamma + b_i + \omega T_{ij} + \varepsilon_{ij}, \quad (5)$$

229 where b_i represents the subject-specific random intercept and ω represents the time effect on the
230 outcome. Note that PLSI mixed-effects model can accommodate additional random effects and
231 other model specifications of fixed effects and interactions, and the model specification for a
232 PLSI mixed-effects model with a random slope was provided in [Additional file 1: Table S1](#). The
233 index coefficient β_l can be explained as the relative direction and importance of X_{ijl} as X_{ijl} is
234 perturbed when scale of β is set and $g(\cdot)$ and other variables are held fixed, and $g(\cdot)$ represents
235 the overall effect of the single index with the mean of longitudinal outcome.

236 Simulation settings

237 Since the NHANES survey dataset does not have time-to-event outcome nor longitudinal
238 outcome, we conducted simulations to demonstrate the PLSI PH model and PLSI mixed-effects
239 model. The coefficients for the 8 environmental factors and three confounding variables were set
240 based on the results from the PLSI linear regression for continuous triglycerides. We kept the
241 original direction of these associations and the absolute rank for each environment factor, and set
242 the effect sizes in a wider range to be more distinguishable (see details in [Table 3](#) and [Table 4](#)).
243 Moreover, we considered the link function $g(\cdot)$ to be either $g(x) = x$ to facilitate the direct
244 comparison with the parametric models, or as a quadratic function $g(x) = x^2$ to mimic the
245 scenario with nonlinear effects and pair-wise interactions between the exposures as
246 $g(\sum_{j=1}^8 \beta_j X_j) = \beta_1^2 X_1^2 + \dots + \beta_8^2 X_8^2 + 2\beta_1\beta_2 X_1 X_2 + \dots + 2\beta_7\beta_8 X_7 X_8$, or a more complex
247 function $g(x) = 0.2x^3 - x^2 + 3x$ to demonstrate higher-order nonlinear effects and
248 interactions, such as three-way interactions. Furthermore, we visualized the interaction effects of
249 two variables by plotting the stratified effect of one variable when fixing the other variables at
250 various levels. Time-to-event outcomes were generated using model (4) with $\lambda_0 = 1$ in the
251 identity link function scenario, $\lambda_0 = 1/\exp(2)$ in the quadratic link function scenario and in the
252 cubic polynomial link function scenario; with a censoring rate as 20% in all of them.

253 Longitudinal outcomes were generated using model (5) with t_{ij} ranged [1, 6] and $\omega = 1$. The
254 number of possible observations for each subject was assumed to vary randomly between 2 and
255 6. The errors followed a first order autoregressive process (i.e. AR(1)), with the autocorrelation
256 as 0.4 and standard deviation as 1.5 to mimic decreasing dependence with time. All details of
257 data generation used in these simulations are included in the R markdown file ([Additional file 3](#)).

258 Performance evaluation

259 In all analyses, the estimated coefficients for the 8 environmental factors and confounders were
260 reported. Ranks based on the absolute values of estimated coefficients were presented to evaluate
261 the relative importance of each environmental factor, and squares of estimated coefficients were
262 shown to represent the respective proportion of contribution to the single index. For all models,
263 the standard errors of coefficient estimates and of the estimated link function were estimated
264 using 500 runs of bootstrapping samples and used to construct the 95% confidence intervals
265 (CIs). We compared the performance of each PLSI model with its counterpart parametric model.
266 The estimated coefficients of 8 environmental factors from the parametric counterpart models
267 were reported in both original values and scaled values to have L_2 norm of 1 for comparison.

268 **Statistical software**

269 All statistical analyses were performed using statistical software R 3.5.0. R codes for the PLSI
270 models for different types of outcomes were developed using ‘gam’, ‘qgam’ or ‘gamm’ function
271 call from ‘mgcv’ or ‘qgam’ package. Linear regression and logistic regression were fit using
272 ‘glm’ function, and quantile regressions using ‘rq’ function in the ‘quantreg’ package. Cox PH
273 model was fitted using ‘coxph’ function from ‘survival’ package, and linear mixed-effects model
274 using ‘lme’ function from ‘nlme’ package. All descriptive and analytical codes were provided as
275 an R Markdown document in [Additional file 3](#).

276 **Results**

277 **Continuous triglycerides: PLSI mean regression**

278 We applied the PLSI linear regression and multivariable linear regression to study the
279 associations of the 8 environmental factors with continuous triglycerides, and summarized the
280 estimates in [Figure 2](#) (numerical results in [Additional file 1: Table S2](#)). The ranks, estimated
281 coefficients, and directions were similar between these two models, and the estimated link

282 function was close to be linear (Additional file 1: Figure S3). As the estimated link function was
283 monotone and increasing, the positive estimates indicated a positive association with
284 triglycerides. Specifically, a-Tocopherol had a $\hat{\beta}_1 = 0.612$ and 95% CI of (0.517, 0.707),
285 indicating that a-Tocopherol had the strongest positive association with triglycerides among the 8
286 factors, and made about 37.4% contribution to the single index; trans-b-carotene had the most
287 negative association of $\hat{\beta}_8 = -0.383$. These results were consistent with original results from
288 Patel's study, which also observed a-Tocopherol with the strongest positive and trans-b-carotene
289 with the strongest negative association with triglycerides (48). As the 8 environmental factors
290 showed both positive and negative associations with triglycerides, this application highlighted
291 the need of statistical methods to accommodate both directional effects for studying multiple
292 environmental exposures. Sensitivity analysis including all 22 environmental factors (Additional
293 file 1: Table S3) showed that the conclusions on the important environmental factors were
294 consistent. The 8 selected environmental factors consistently showed top ranks among the 22
295 factors, except for PCB194 which was highly correlated with other PCBs. When there are many
296 highly correlated exposures ($r > 0.9$), we also recommend using p-values to rank the importance
297 of variables in addition to the absolute coefficient values, which can be inflated by
298 multicollinearity (65).

299 Continuous triglycerides: PLSI quantile regression

300 We applied the PLSI quantile regression to study the associations between 8 exposures and
301 three quartiles (25th, 50th, and 75th percentiles) of triglycerides and summarized the main results
302 in Figure 3 (numerical results in Additional file 1: Table S4). We observed that the estimated
303 link functions for all three quartiles were increasing and close to linear (Additional file 1: Figure
304 S4), which explained the similarities between the results of the PLSI quantile regressions and

305 regular quantile regressions. In addition, the 8 environmental factors showed fairly consistent
306 associations across the three quartiles of triglycerides. For example, a-Tocopherol was the factor
307 having the strongest positive association with triglycerides and trans-b-carotene was the factor
308 having the strongest negative association with triglycerides at all three quartiles.

309 Binary triglycerides: PLSI logistic regression

310 For dichotomized triglycerides, the ranks and estimates from PLSI logistic regression and
311 multivariable logistic regression are shown in [Figure 4](#) (numerical results in [Additional file 1:
312 Table S5](#)), which demonstrated similar results from these two models. The estimated link
313 function by PLSI logistic regression was monotone increasing and close to be linear ([Additional
314 file 1: Figure S5](#)). Thus, the estimated directions can be interpreted qualitatively and the
315 estimated coefficients represented the relative importance of each exposure on the log odds of
316 high triglycerides. For example, the estimated coefficient of a-Tocopherol was $\hat{\beta}_1 = 0.584$ (95%
317 CI: 0.433-0.735), which represented that a-Tocopherol had the strongest positive association
318 with the odds of high triglycerides among the 8 factors.

319 Simulated time-to-event outcome: PLSI PH model

320 We summarize the simulation results from both PLSI PH model and Cox PH model in [Table
321 3](#). Under the identity link function setting, results from the PLSI PH model and the conventional
322 Cox PH model were very similar as expected, and both close to the true values. The PLSI PH
323 model estimated the link function to be close to the true linear function ([Additional file 1: Figure
324 S6 \(a\)](#)). Under the quadratic link function setting, results from the PLSI PH model were still
325 consistent to true coefficients, but the conventional Cox PH model failed for most of the
326 environmental factors because the linear model assumption was insufficient. The PLSI PH model
327 also captured the U-shape and estimated the link function close to the true quadratic function

328 ([Additional file 1: Figure S6 \(b\)](#)). Stratified plots ([Additional file 1: Figure S7](#)) showed that a-
329 Tocopherol had different effects on the outcome when trans-b-carotene was set at its 10th, 50th,
330 and 90th percentiles, indicating the existence of an interaction between a-Tocopherol and trans-b-
331 carotene in this scenario. Results for cubic polynomial link function ([Additional file 1: Table](#)
332 [S6/Figure S8](#)) presented good performance in coefficient and link function estimations,
333 suggesting that PLSI models are able to handle complex higher-order interactions among
334 environmental factors.

335 ([Table 3](#) should appear here)

336 Simulated longitudinal outcome: PLSI mixed-effects model

337 The results from PLSI mixed-effects model and linear mixed-effects model under identify or
338 quadratic link function are presented in [Table 4](#). Under the identity link function setting, the
339 PLSI mixed-effects model estimated all coefficients close to the true coefficients with correct
340 directions, and conventional linear mixed-effects model also had similar estimations. The
341 estimated link function by PLSI mixed-effects model was close to the true linear function
342 ([Additional file 1: Figure S9 \(a\)](#)). Under the quadratic link function setting, the results from PLSI
343 mixed-effects model were still consistent; however, the conventional linear mixed-effects model
344 clearly showed biased results for some factors like PCB194. The estimated link function by PLSI
345 mixed-effects model had a U-shape and was close to the true quadratic function ([Additional file](#)
346 [1: Figure S9 \(b\)](#)).

347 ([Table 4](#) should appear here)

348 Discussion

349 We presented five PLSI models aiming to provide a unified family of statistical models to assess
350 the joint effects of environmental exposures on four types of health outcomes: continuous,

351 categorical, time-to-event, and longitudinal outcomes. We demonstrated the flexibility and
352 effectiveness of this PLSI family for modeling various types of outcomes using NHANES data
353 supplemented with simulations. One contribution of this work is that the novel modeling options
354 under the PLSI framework complement existing methods and address some common statistical
355 challenges in the analysis of multiple environmental exposures, such as mixed directions,
356 interactions, and non-linear effects. Another contribution is that coherent computation algorithms
357 are developed for all the PLSI models and implemented using the existing R packages, which
358 can facilitate direct applications in practice and reproducible research.

359 In our analyses of the cross-sectional NHANES studies for continuous and binary
360 triglycerides by PLSI models, we found that the 8 environmental factors exhibited mixed
361 directional associations with the outcome, with α -Tocopherol having the strongest positive
362 association and trans-b-carotene having the strongest negative association with triglycerides. α -
363 Tocopherol and carotenes are transported in serum with HDL and LDL, and the level of serum α -
364 Tocopherol depends on serum lipids (66, 67). The strong positive association between α -
365 Tocopherol and triglycerides is expected (48), and the negative association between b-carotene
366 and triglycerides is supported by previous studies (68, 69). Our results were consistent with the
367 results of previously known and validated environmental chemical factors correlated with
368 triglycerides (48), clearly demonstrating the value of PLSI models as a flexible and useful tool
369 for analyzing complex exposures. Using additional simulations for time-to-event and
370 longitudinal outcomes, we showed that the PLSI models could correctly identify the directions
371 and magnitudes of associations for these environmental factors in scenarios with different types
372 of outcomes.

373 In our NHANES applications of studying triglycerides continuously and categorically, we
374 estimated that the link functions of PLSI models were very close to be linear, which were also

375 reflected by the similar results with their counterpart parametric models. In general, standard
376 errors from the PLSI models were larger than those from their counterpart parametric models,
377 which was expected as the former are semiparametric models.

378 We also conducted another sensitivity weighted analysis incorporating the laboratory
379 subsample C weights from NHANES 2003-2004 cycle (following general guideline to use the
380 weights from “least common denominator”) (70), and the weighted results (Additional file 1:
381 Table S7) were similar with the results from unweighted models. Note that most of the PLSI
382 models are readily incorporate weights in R function codes (Additional file 3).

383 Interaction among multiple correlated environmental factors is very common, and it has
384 been long appreciated that the co-exposures may have synergistic (additive or multiplicative) or
385 antagonistic effects on health outcomes (71). For parametric models, it’s difficult to directly
386 model the interaction effects among co-exposures if we don’t know the ‘degree of interaction’.
387 However, PLSI models can handle the interaction easily through the unknown link function as
388 we evaluated using the simulations. Specifically, in our simulated time-to-event and longitudinal
389 analyses with quadratic link function, which reflected both the pairwise interactions and non-
390 linear quadratic effects, both PLSI PH model and PLSI mixed-effects were able to capture the U-
391 shape link function and correct direction and importance of the environmental factors, while
392 parametric models failed in most factors because the parametric assumptions were no longer
393 satisfied. For more complex (higher-order) interactions, the flexibility of the nonparametric link
394 function can incorporate the effects of these interactions (72). Therefore, PLSI models readily
395 accommodate the factors showing non-linear or interactive effect on the health outcome.

396 There are other ways and models using various definitions of weighted sums to model the
397 joint effect for multiple environmental components. For example, molar sums were used to show
398 relationships between prenatal phenol and phthalate exposures and birth outcome (73), and a

399 potency-weighted sum was used to calculate phthalates exposures among reproductive-aged
400 women (74). The weights for environmental factors can be calculated from their expected
401 potency relative to a reference factor, like the common cases in toxicology (75), or based on their
402 percent contribution to the total mixture effect, like WQS (9). PLSI models can be considered as
403 one of these weighting approaches, and their advantages from the semiparametric structure are
404 evident compared with existing methods, especially for the scenarios when the environmental
405 exposures have mixed-directional associations and/or a potential high-degree interaction.
406 Meanwhile, due to the flexibility of the nonparametric link function, PLSI models can represent
407 complex joint effects more than additive structures (76), which is commonly encountered since
408 environmental exposures may act together in a biological sense via a shared mechanistic
409 pathway (4). The ability of handling various types of outcomes is another important advantage of
410 the proposed PLSI framework. This is important because, with the accumulation of
411 environmental exposure measurements and development of data collection methods, time-to-
412 event or longitudinal studies are desired to explore the associations over time.

413 In this study, the coherent algorithms for PLSI models are based on the ‘gam’ and ‘gamm’
414 functions from ‘mgcv’ package and ‘qgam’ function from ‘qgam’ package in R, which includes
415 many of the generalized additive model (GAM) fitting techniques developed by Simon Wood et
416 al (77). The rationale behind the algorithms is to use ‘gam’, ‘qgam’ or ‘gamm’ call (usually
417 using penalized regression splines or similar smoothers) to profile out the smooth model
418 coefficients and smoothing parameters for estimation of the link function contained in PLSI
419 model, leaving only a finite parameter vector to be estimated by a general purpose optimizer.
420 Based on this algorithm, it is easy to adapt the models to include multiple single index terms,
421 parametric terms, and further smoothing. We have compared the estimates for single index
422 models among different iterative procedures using existing packages (e.g., projection pursuit

423 regression with one term using ‘ppr’ function; ‘sim.est’ function from ‘simest’ package) in
424 various simulations, and they have similar estimation performance. We finally chose ‘gam’ call
425 series because of its flexibility for covariate adjustment and ability of modeling various types of
426 outcomes. This ‘gam’, ‘qgam’, ‘gamm’ call approach has demonstrated efficient and robust
427 performance in our numerical studies, and we believe this coherent algorithm strategy wrapped
428 as a toolbox is beneficial for practical application.

429 The PLSI models considered here may not be directly applicable to extreme high-
430 dimensional settings, for which we could consider using extensions with adaptive LASSO (78),
431 smoothly clipped absolute deviation penalty (79), and smooth-threshold estimating equations
432 (80). Another future research direction is to extend from the single index to multiple-index
433 models, such as the projection pursuit regression (81), so that more complex data structures and
434 exposure effect patterns can be captured and modeled.

435 **Conclusions**

436 A family of PLSI models exemplified great value of identifying important components among
437 environmental exposures when they demonstrate associations in various directions and complex
438 non-linear relationships between the exposures and outcome.

439 **Additional files**

440 **Addition file 1: Figure S1.** Data flow diagram for deriving 800 subjects and 8 environmental
441 factors. **Figure S2.** Correlation matrix of Pearson correlation coefficient of 22 factors and
442 triglycerides in NHANES 2002-2003 (N=800). **Table S1.** PLSI generalized linear regression for
443 ordinal, multinomial, and count outcomes and PLSI mixed-effects model with random slope for
444 longitudinal outcome. **Table S2.** Results from PLSI linear regression and multivariable linear

445 regression in NHANES 2002-2003. **Figure S3.** Estimated link function by PLSI linear regression in
446 NHANES 2002-2003. **Table S3.** Sensitivity analysis results from PLSI linear regression and
447 multivariable linear regression in NHANES 2002-2003 with 22 environmental factors. **Tables S4.1-**
448 **S4.3.** Results from PLSI quantile regressions and multivariable quantile regression at three
449 quantiles (25th, 50th, and 75th percentiles) of triglycerides in NHANES 2002-2003. **Figure S4.**
450 Estimated link functions by PLSI quantile regressions at three quartiles in NHANES 2002-2003. (a)
451 25th percentile; (b) 50th percentile; (c) 75th percentile. **Table S5.** Results from PLSI logistic
452 regression and multivariable logistic regression in NHANES 2002. **Figure S5.** Estimated link
453 function by PLSI logistic regression in NHANES 2002-2003. **Figure S6.** Estimated link functions by
454 PLSI PH model in simulated time-to-event study. (a) identity link function; (b) quadratic link
455 function. **Figure S7.** Stratified effect of a-Tocopherol with 95% confidence intervals when the
456 variable of trans-b-carotene fixed at 10%, 50%, and 90% percentile and other factors fixed as
457 median values. **Table S6.** Simulation results from PLSI PH model and Cox PH model for link
458 function $g(x) = 0.2x^3 - x^2 + 3x$. **Figure S8.** Estimated link functions by PLSI PH model in
459 simulated time-to-event study with link function $g(x) = 0.2x^3 - x^2 + 3x$. **Figure S9.** Estimated
460 link functions by PLSI mixed-effects model in simulated longitudinal study. (a) identity link
461 function; (b) quadratic link function. **Table S7.** Sensitivity analysis results from weighted PLSI
462 linear regression and weighted linear regression in NHANES 2002-2003 using NHANES laboratory
463 subsample C weights.

464 **Addition file 2:** cleaning dataset of 800 subjects from NHANES 2003-2004 cycle. Variables
465 include respondent sequence number of subject, outcome triglyceride, 22 environmental
466 factors, 3 demographic confounding variables, and laboratory subsample C weight.

467 **Addition file 3:** R markdown document demonstrating all descriptive and analytical process of
468 this article.

469 **Abbreviations**

470 AR: autoregressive process; BKMR: Bayesian kernel machine regression; NHANES: National
471 Health and Nutrition Examination Survey; NIEHS: National Institute of Environmental Health
472 Sciences; PH: proportional hazards; PLSI: partial-linear single-index; PIP: posterior inclusion
473 probability; QR: quantile regression; WQS: weighted quantile sum regression

474 **Declarations**

475 **Ethics approval and consent to participate**

476 Subjects provided the written informed consent, and the institutional review board of the
477 National Center for Health Statistics approved the survey for NHANES study.

478 **Consent for publications**

479 Not applicable.

480 **Availability of data and materials**

481 The dataset used and/or analyzed during the current study supporting the conclusions of this
482 article is included within the additional file.

483 **Competing interests**

484 The authors declare that they have no competing interests.

485 **Funding**

486 This work is partially supported by UG3/UH3OD023305 and 4P30ES000260-52 from the National
487 Institutes of Health.

488 **Authors' contributions**

489 YWang and MLiu: Performed data curation, conducted statistical analyses and prepared original
490 manuscript draft. YWang, YWu, MLee, and PJ: Designed the algorithm and performed
491 simulations. MJ, LT and MLiu: Directed the data set collection and quality control, acquired
492 funding to support this analysis, contributed to literature review and reviewed the manuscript.
493 All authors read and approve the final manuscript.

494 **Acknowledgements**

495 The contributions of the subjects in the NHANES study are gratefully acknowledged.

496 **References**

- 498 1. Wild CP. Complementing the genome with an "exposome": The outstanding challenge of
499 environmental exposure measurement in molecular epidemiology. *Cancer Epidem Biomar.*
500 2005;14(8):1847-50.
- 501 2. Stafoggia M, Breitner S, Hampel R, Basagana X. Statistical Approaches to Address
502 Multi-Pollutant Mixtures and Multiple Exposures: the State of the Science. *Curr Environ Health*
503 *Rep.* 2017;4(4):481-90.
- 504 3. Sanders AP, Claus Henn B, Wright RO. Perinatal and Childhood Exposure to Cadmium,
505 Manganese, and Metal Mixtures and Effects on Cognition and Behavior: A Review of Recent
506 Literature. *Curr Environ Health Rep.* 2015;2(3):284-94.
- 507 4. Hamra GB, Buckley JP. Environmental exposure mixtures: questions and methods to
508 address them. *Curr Epidemiol Rep.* 2018;5(2):160-5.
- 509 5. NIEHS Strategic Plan 2018–2023 2018 [Available from:
510 [https://www.niehs.nih.gov/about/strategicplan/index.cfm#:~:text=The%20NIEHS%20strategic%
511 20plan%202018,EHS%20Through%20Stewardship%20and%20Support](https://www.niehs.nih.gov/about/strategicplan/index.cfm#:~:text=The%20NIEHS%20strategic%20plan%202018,EHS%20Through%20Stewardship%20and%20Support).
- 512 6. Billionnet C, Sherrill D, Annesi-Maesano I, Study G. Estimating the Health Effects of
513 Exposure to Multi-Pollutant Mixture. *Ann Epidemiol.* 2012;22(2):126-41.
- 514 7. Mann RM, Hyne RV, Choung CB, Wilson SP. Amphibians and agricultural chemicals:
515 review of the risks in a complex environment. *Environ Pollut.* 2009;157(11):2903-27.
- 516 8. Chaumont A, Nickmilder M, Dumont X, Lundh T, Skerfving S, Bernard A. Associations
517 between proteins and heavy metals in urine at low environmental exposures: Evidence of reverse
518 causality. *Toxicol Lett.* 2012;210(3):345-52.

- 519 9. Carrico C, Gennings C, Wheeler DC, Factor-Litvak P. Characterization of Weighted
520 Quantile Sum Regression for Highly Correlated Data in a Risk Analysis Setting. *J Agr Biol*
521 *Envir St.* 2015;20(1):100-20.
- 522 10. Czarnota J, Gennings C, Colt JS, De Roos AJ, Cerhan JR, Severson RK, et al. Analysis
523 of Environmental Chemical Mixtures and Non-Hodgkin Lymphoma Risk in the NCI-SEER NHL
524 Study. *Environ Health Persp.* 2015;123(10):965-70.
- 525 11. Bobb JF, Valeri L, Claus Henn B, Christiani DC, Wright RO, Mazumdar M, et al.
526 Bayesian kernel machine regression for estimating the health effects of multi-pollutant mixtures.
527 *Biostatistics.* 2015;16(3):493-508.
- 528 12. Valeri L, Mazumdar MM, Bobb JF, Henn BC, Rodrigues E, Sharif OIA, et al. The Joint
529 Effect of Prenatal Exposure to Metal Mixtures on Neurodevelopmental Outcomes at 20-40
530 Months of Age: Evidence from Rural Bangladesh. *Environ Health Persp.* 2017;125(6).
- 531 13. Zhang YQ, Dong TY, Hu WY, Wang X, Xu B, Lin ZN, et al. Association between
532 exposure to a mixture of phenols, pesticides, and phthalates and obesity: Comparison of three
533 statistical models. *Environ Int.* 2019;123:325-36.
- 534 14. Keil AP, Buckley JP, O'Brien KM, Ferguson KK, Zhao S, White AJ. A Quantile-Based
535 g-Computation Approach to Addressing the Effects of Exposure Mixtures. *Environ Health*
536 *Perspect.* 2020;128(4):47004.
- 537 15. Levin-Schwartz Y, Gennings C, Schnaas L, Del Carmen Hernandez Chavez M, Bellinger
538 DC, Tellez-Rojo MM, et al. Time-varying associations between prenatal metal mixtures and
539 rapid visual processing in children. *Environ Health.* 2019;18(1):92.
- 540 16. Zhang L, Kim I. Semiparametric Bayesian kernel survival model for evaluating pathway
541 effects. *Stat Methods Med Res.* 2019;28(10-11):3301-17.
- 542 17. Gibson EA, Nunez Y, Abuawad A, Zota AR, Renzetti S, Devick KL, et al. An overview
543 of methods to address distinct research questions on environmental mixtures: an application to
544 persistent organic pollutants and leukocyte telomere length. *Environ Health-Glob.* 2019;18(1).
- 545 18. Ichimura H. Semiparametric Least-Squares (Sls) and Weighted Sls Estimation of Single-
546 Index Models. *J Econometrics.* 1993;58(1-2):71-120.
- 547 19. Horowitz JL, Hardle W. Direct semiparametric estimation of single-index models with
548 discrete covariates. *J Am Stat Assoc.* 1996;91(436):1632-40.
- 549 20. Wang JL, Xue LG, Zhu LX, Chong YS. Estimation for a Partial-Linear Single-Index
550 Model. *Ann Stat.* 2010;38(1):246-74.
- 551 21. Hardle W, Hall P, Ichimura H. Optimal Smoothing in Single-Index Models. *Ann Stat.*
552 1993;21(1):157-78.
- 553 22. Carroll RJ, Fan JQ, Gijbels I, Wand MP. Generalized partially linear single-index
554 models. *J Am Stat Assoc.* 1997;92(438):477-89.
- 555 23. Yi GY, He WQ, Liang H. Analysis of correlated binary data under partially linear single-
556 index logistic models. *J Multivariate Anal.* 2009;100(2):278-90.
- 557 24. Wang W. Proportional hazards regression models with unknown link function and time-
558 dependent covariates. *Stat Sinica.* 2004;14(3):885-905.
- 559 25. Huang JHZ, Liu LX. Polynomial spline estimation and inference of proportional hazards
560 regression models with flexible relative risk form. *Biometrics.* 2006;62(3):793-802.
- 561 26. Sun J, Kopciuk KA, Lu XW. Polynomial spline estimation of partially linear single-index
562 proportional hazards regression models. *Comput Stat Data An.* 2008;53(1):176-88.
- 563 27. Li JB, Zhang RQ. Partially varying coefficient single index proportional hazards
564 regression models. *Comput Stat Data An.* 2011;55(1):389-400.

- 565 28. Bai Y, Fung WK, Zhu ZY. Penalized quadratic inference functions for single-index
566 models with longitudinal data. *J Multivariate Anal.* 2009;100(1):152-61.
- 567 29. Li GR, Zhu LX, Xue LG, Feng SY. Empirical likelihood inference in partially linear
568 single-index models for longitudinal data. *J Multivariate Anal.* 2010;101(3):718-32.
- 569 30. Xu PR, Zhu LX. Estimation for a marginal generalized single-index longitudinal model. *J*
570 *Multivariate Anal.* 2012;105(1):285-99.
- 571 31. Zhao WH, Lian H, Liang H. GEE analysis for longitudinal single-index quantile
572 regression. *J Stat Plan Infer.* 2017;187:78-102.
- 573 32. Stoker TM. Consistent Estimation of Scaled Coefficients. *Econometrica.*
574 1986;54(6):1461-81.
- 575 33. Hardle W, Stoker TM. Investigating Smooth Multiple-Regression by the Method of
576 Average Derivatives. *J Am Stat Assoc.* 1989;84(408):986-95.
- 577 34. Hardle W, Tsybakov AB. How Sensitive Are Average Derivatives. *J Econometrics.*
578 1993;58(1-2):31-48.
- 579 35. Hristache M, Juditsky A, Spokoiny V. Direct estimation of the index coefficient in a
580 single-index model. *Ann Stat.* 2001;29(3):595-623.
- 581 36. Yu Y, Ruppert D. Penalized spline estimation for partially linear single-index models. *J*
582 *Am Stat Assoc.* 2002;97(460):1042-54.
- 583 37. Xia YC, Hardle W. Semi-parametric estimation of partially linear single-index models. *J*
584 *Multivariate Anal.* 2006;97(5):1162-84.
- 585 38. Liang H, Liu X, Li RZ, Tsai CL. Estimation and Testing for Partially Linear Single-Index
586 Models. *Ann Stat.* 2010;38(6):3811-36.
- 587 39. Chaudhuri P. Global Nonparametric-Estimation of Conditional Quantile Functions and
588 Their Derivatives. *J Multivariate Anal.* 1991;39(2):246-69.
- 589 40. Chaudhuri P, Doksum K, Samarov A. On average derivative quantile regression. *Ann*
590 *Stat.* 1997;25(2):715-44.
- 591 41. Wu TZ, Yu KM, Yu Y. Single-index quantile regression. *J Multivariate Anal.*
592 2010;101(7):1607-21.
- 593 42. Kong EF, Xia YC. A Single-Index Quantile Regression Model and Its Estimation.
594 *Economet Theor.* 2012;28(4):730-68.
- 595 43. Lv YZ, Zhang RQ, Zhao WH, Liu JC. Quantile regression and variable selection of
596 partial linear single-index model. *Ann I Stat Math.* 2015;67(2):375-409.
- 597 44. Ma SJ, He XM. Inference for Single-Index Quantile Regression Models with Profile
598 Optimization. *Ann Stat.* 2016;44(3):1234-68.
- 599 45. Lai P, Li GR, Lian H. Quadratic inference functions for partially linear single-index
600 models with longitudinal data. *J Multivariate Anal.* 2013;118:115-27.
- 601 46. Li GR, Lai P, Lian H. Variable selection and estimation for partially linear single-index
602 models with longitudinal data. *Stat Comput.* 2015;25(3):579-93.
- 603 47. Li JB, Lian H, Jiang XJ, Song XY. Estimation and testing for time-varying quantile
604 single-index models with longitudinal data. *Comput Stat Data An.* 2018;118:66-83.
- 605 48. Patel CJ, Cullen MR, Ioannidis JPA, Butte AJ. Systematic evaluation of environmental
606 factors: persistent pollutants and nutrients correlated with serum lipid levels. *Int J Epidemiol.*
607 2012;41(3):828-43.
- 608 49. Zipf G, Chiappa M, Porter KS, Ostchega Y, Lewis BG, Dostal J. National health and
609 nutrition examination survey: plan and operations, 1999-2010. *Vital Health Stat 1.* 2013(56):1-
610 37.
- 611 50. Weisberg S, Welsh AH. Adapting for the Missing Link. *Ann Stat.* 1994;22(4):1674-700.

- 612 51. Di Angelantonio E, Sarwar N, Perry P, Kaptoge S, Ray KK, Thompson A, et al. Major
613 Lipids, Apolipoproteins, and Risk of Vascular Disease. *Jama-J Am Med Assoc.*
614 2009;302(18):1993-2000.
- 615 52. Bind MA, Peters A, Koutrakis P, Coull B, Vokonas P, Schwartz J. Quantile Regression
616 Analysis of the Distributional Effects of Air Pollution on Blood Pressure, Heart Rate Variability,
617 Blood Lipids, and Biomarkers of Inflammation in Elderly American Men: The Normative Aging
618 Study. *Environ Health Persp.* 2016;124(8):1189-98.
- 619 53. Burgette LF, Reiter JP, Miranda ML. Exploratory Quantile Regression With Many
620 Covariates An Application to Adverse Birth Outcomes. *Epidemiology.* 2011;22(6):859-66.
- 621 54. Ratcliff R, Thapar A, McKoon G. Individual differences, aging, and IQ in two-choice
622 tasks. *Cognitive Psychol.* 2010;60(3):127-57.
- 623 55. Jung SH. Quasi-likelihood for median regression models. *J Am Stat Assoc.*
624 1996;91(433):251-7.
- 625 56. Koenker R, Bassett G. Regression Quantiles. *Econometrica.* 1978;46(1):33-50.
- 626 57. Koenker R, Hallock KF. Quantile regression. *J Econ Perspect.* 2001;15(4):143-56.
- 627 58. Wei Y, Pere A, Koenker R, He XM. Quantile regression methods for reference growth
628 charts. *Stat Med.* 2006;25(8):1369-82.
- 629 59. Expert Panel on Detection E, Treatment of High Blood Cholesterol in A. Executive
630 Summary of The Third Report of The National Cholesterol Education Program (NCEP) Expert
631 Panel on Detection, Evaluation, And Treatment of High Blood Cholesterol In Adults (Adult
632 Treatment Panel III). *JAMA.* 2001;285(19):2486-97.
- 633 60. Cox DR. Regression Models and Life-Tables. *J R Stat Soc B.* 1972;34(2):187-+.
- 634 61. Cox DR. Partial Likelihood. *Biometrika.* 1975;62(2):269-76.
- 635 62. Dempster AP, Laird NM, Rubin DB. Maximum Likelihood from Incomplete Data Via
636 Em Algorithm. *J Roy Stat Soc B Met.* 1977;39(1):1-38.
- 637 63. Laird NM, Ware JH. Random-Effects Models for Longitudinal Data. *Biometrics.*
638 1982;38(4):963-74.
- 639 64. Rubin DB. Inference and Missing Data. *Biometrika.* 1976;63(3):581-90.
- 640 65. Wold S, Ruhe A, Wold H, Dunn WJ. The Collinearity Problem in Linear-Regression -
641 the Partial Least-Squares (Pls) Approach to Generalized Inverses. *Siam J Sci Stat Comp.*
642 1984;5(3):735-43.
- 643 66. Ogihara T, Miki M, Kitagawa M, Mino M. Distribution of Tocopherol among Human-
644 Plasma Lipoproteins. *Clin Chim Acta.* 1988;174(3):299-305.
- 645 67. Winbauer AN, Pingree SS, Nuttall KL. Evaluating serum alpha-tocopherol (vitamin E) in
646 terms of a lipid ratio. *Ann Clin Lab Sci.* 1999;29(3):185-91.
- 647 68. Vanvliet T, Schreurs WHP, Vandenberg H. Intestinal Beta-Carotene Absorption and
648 Cleavage in Men - Response of Beta-Carotene and Retinyl Esters in the Triglyceride-Rich
649 Lipoprotein Fraction after a Single Oral Dose of Beta-Carotene. *Am J Clin Nutr.*
650 1995;62(1):110-6.
- 651 69. Redlich CA, Chung JS, Cullen MR, Blaner WS, Van Bennekum AM, Berglund L. Effect
652 of long-term beta-carotene and vitamin A on serum cholesterol and triglyceride levels among
653 participants in the Carotene and Retinol Efficacy trial (CARET) (vol 143, pg 427, 1999).
654 *Atherosclerosis.* 1999;145(2):423-+.
- 655 70. Johnson CL, Paulose-Ram R, Ogden CL, Carroll MD, Kruszon-Moran D, Dohrmann
656 SM, et al. National health and nutrition examination survey: analytic guidelines, 1999-2010.
657 *Vital Health Stat 2.* 2013(161):1-24.

658 71. Walter SD, Holford TR. Additive, Multiplicative, and Other Models for Disease Risks.
659 Am J Epidemiol. 1978;108(5):341-6.
660 72. Radchenko P. High dimensional single index models. J Multivariate Anal. 2015;139:266-
661 82.
662 73. Wolff MS, Engel SM, Berkowitz GS, Ye X, Silva MJ, Zhu C, et al. Prenatal phenol and
663 phthalate exposures and birth outcomes. Environ Health Perspect. 2008;116(8):1092-7.
664 74. Varshavsky JR, Zota AR, Woodruff TJ. A Novel Method for Calculating Potency-
665 Weighted Cumulative Phthalates Exposure with Implications for Identifying Racial/Ethnic
666 Disparities among U.S. Reproductive-Aged Women in NHANES 2001-2012. Environ Sci
667 Technol. 2016;50(19):10616-24.
668 75. Howard GJ, Webster TF. Contrasting theories of interaction in epidemiology and
669 toxicology. Environ Health Perspect. 2013;121(1):1-6.
670 76. VanderWeele TJ. On the Distinction Between Interaction and Effect Modification.
671 Epidemiology. 2009;20(6):863-71.
672 77. Pedersen EJ, Miller DL, Simpson GL, Ross N. Hierarchical generalized additive models
673 in ecology: an introduction with mgcv. PeerJ. 2019;7:e6876.
674 78. Foster JC, Taylor JMG, Nan B. Variable selection in monotone single-index models via
675 the adaptive LASSO. Stat Med. 2013;32(22):3944-54.
676 79. Yang H, Yang J. A robust and efficient estimation and variable selection method for
677 partially linear single-index models. J Multivariate Anal. 2014;129:227-42.
678 80. Lai P, Wang QH, Lian H. Bias-corrected GEE estimation and smooth-threshold GEE
679 variable selection for single-index models with clustered data. J Multivariate Anal.
680 2012;105(1):422-32.
681 81. Friedman JH, Stuetzle W. Projection Pursuit Regression. J Am Stat Assoc.
682 1981;76(376):817-23.

683

684

685

686

687

688

689

690

691

692

693

694 **Table 3** Simulation results from PLSI PH model and Cox PH model

Variable	True rank	True coefficient	PLSI PH rank	PLSI PH estimate	PLSI PH 95% CI	PLSI PH Proportion of contribution (%)	Cox PH rank	Cox PH original estimate	Cox PH original 95% CI	Cox PH normed estimate	Cox PH normed 95% CI
Identity link function											
Environmental factors											
a-Tocopherol	1	0.560	1	0.546	(0.437, 0.656)	29.9	1	0.558	(0.428, 0.688)	0.546	(0.446, 0.646)
g-tocopherol	2	0.490	2	0.500	(0.427, 0.572)	25.0	2	0.511	(0.417, 0.605)	0.500	(0.428, 0.571)
Retinyl-palmitate	3	0.420	3	0.408	(0.297, 0.520)	16.7	3	0.418	(0.310, 0.526)	0.409	(0.301, 0.516)
Retinol	7	0.140	7	0.122	(0.029, 0.216)	1.5	7	0.125	(0.029, 0.221)	0.122	(0.034, 0.210)
3,3,4,4,5-pncb	8	0.070	8	0.059	(-0.039, 0.158)	0.4	8	0.061	(-0.040, 0.161)	0.059	(-0.033, 0.151)
PCB194	6	-0.210	6	-0.207	(-0.346, -0.068)	4.3	6	-0.212	(-0.351, -0.074)	-0.208	(-0.329, -0.087)
2.3.4.6.7.8.hxcdf	5	-0.280	5	-0.270	(-0.356, -0.183)	7.3	5	-0.275	(-0.367, -0.183)	-0.269	(-0.354, -0.185)
trans.b.carotene	4	-0.350	4	-0.388	(-0.467, -0.310)	15.1	4	-0.397	(-0.493, -0.302)	-0.389	(-0.465, -0.313)
Covariates											
Age		0.005		0.009	(0.001, 0.017)			0.009	(0.002, 0.016)		
Sex (female)		-0.076		-0.039	(-0.216, 0.138)			-0.039	(-0.217, 0.138)		
Ethnicity											
Non-Hispanic white		Ref		Ref				Ref			
Non-Hispanic black		-0.138		-0.135	(-0.367, 0.097)			-0.135	(-0.361, 0.091)		
Mexican American		0.175		0.114	(-0.116, 0.344)			0.114	(-0.107, 0.335)		
Other race		0.409		0.528	(0.118, 0.937)			0.528	(0.077, 0.978)		
Other Hispanic		0.355		0.477	(-0.021, 0.975)			0.477	(0.018, 0.936)		
Quadratic link function											
Environmental factors											
a-Tocopherol	1	0.560	1	0.526	(0.403, 0.648)	27.6	1	0.289	(0.124, 0.455)	0.861	(0.621, 1.101)
g-tocopherol	2	0.490	2	0.513	(0.296, 0.730)	26.3	3	0.098	(-0.011, 0.207)	0.292	(-0.024, 0.607)
Retinyl-palmitate	3	0.420	3	0.445	(0.231, 0.659)	19.8	6	0.037	(-0.088, 0.161)	0.109	(-0.253, 0.470)
Retinol	7	0.140	7	0.161	(0.041, 0.281)	2.6	4	-0.041	(-0.154, 0.072)	-0.122	(-0.465, 0.222)
3,3,4,4,5-pncb	8	0.070	8	0.061	(-0.023, 0.146)	0.4	8	0.013	(-0.102, 0.128)	0.040	(-0.305, 0.384)
PCB194	6	-0.210	6	-0.208	(-0.322, -0.093)	4.3	7	0.020	(-0.132, 0.172)	0.059	(-0.338, 0.457)
2.3.4.6.7.8.hxcdf	5	-0.280	5	-0.252	(-0.392, -0.113)	6.4	5	-0.039	(-0.138, 0.061)	-0.115	(-0.445, 0.215)
trans.b.carotene	4	-0.350	4	-0.355	(-0.477, -0.234)	12.6	2	-0.120	(-0.228, -0.012)	-0.358	(-0.637, -0.079)
Covariates											
Age		0.005		0.003	(-0.002, 0.008)			-0.005	(-0.012, 0.003)		
Sex (female)		-0.076		-0.081	(-0.269, 0.108)			-0.103	(-0.297, 0.092)		
Ethnicity											
Non-Hispanic white		Ref		Ref				Ref			
Non-Hispanic black		-0.138		0.044	(-0.211, 0.299)			0.083	(-0.154, 0.320)		
Mexican American		0.175		0.100	(-0.152, 0.352)			0.125	(-0.118, 0.369)		
Other race		0.409		0.186	(-0.438, 0.811)			-0.189	(-0.722, 0.345)		
Other Hispanic		0.355		0.096	(-0.567, 0.759)			-0.096	(-0.634, 0.442)		

695

696

697

698

699

700

701

702

703 **Table 4** Simulation results from PLSI mixed-effects model and linear mixed-effects model

Variable	True rank	True coefficient	PLSI ME rank	PLSI ME estimate	PLSI ME 95% CI	PLSI ME Proportion of contribution (%)	Linear ME rank	Linear ME original estimate	Linear ME original 95% CI	Linear ME normed estimate	Linear ME normed 95% CI
Identity link function											
Environmental factors											
a-Tocopherol	1	0.560	1	0.584	(0.469, 0.698)	34.1	1	0.590	(0.456, 0.723)	0.580	(0.519, 0.642)
g-tocopherol	2	0.490	2	0.481	(0.396, 0.566)	23.1	2	0.490	(0.401, 0.579)	0.482	(0.439, 0.525)
Retinyl-palmitate	3	0.420	3	0.402	(0.284, 0.520)	16.2	3	0.408	(0.302, 0.513)	0.401	(0.336, 0.467)
Retinol	7	0.140	7	0.091	(-0.025, 0.206)	0.8	7	0.088	(-0.011, 0.186)	0.086	(0.027, 0.145)
3,3,4,4,5-pncb	8	0.070	8	0.054	(-0.067, 0.175)	0.3	8	0.058	(-0.047, 0.164)	0.057	(0.000, 0.114)
PCB194	6	-0.210	6	-0.225	(-0.378, -0.072)	5.1	6	-0.236	(-0.372, -0.099)	-0.232	(-0.303, -0.160)
2.3.4.6.7.8.hxcdf	5	-0.280	5	-0.236	(-0.344, -0.128)	5.6	5	-0.241	(-0.331, -0.151)	-0.237	(-0.295, -0.179)
trans.b.carotene	4	-0.350	4	-0.386	(-0.475, -0.297)	14.9	4	-0.392	(-0.486, -0.298)	-0.386	(-0.433, -0.339)
Covariates											
Intercept		0.000		-0.069	(-0.426, 0.287)			-0.074	(-0.486, 0.339)		
Age		0.005		0.011	(0.004, 0.019)			0.011	(0.005, 0.018)		
Sex (female)		-0.076		-0.121	(-0.245, 0.003)			-0.125	(-0.302, 0.051)		
Ethnicity											
Non-Hispanic white		Ref		Ref				Ref			
Non-Hispanic black		-0.138		-0.225	(-0.368, -0.082)			-0.231	(-0.450, -0.012)		
Mexican American		0.175		0.030	(-0.123, 0.184)			0.027	(-0.195, 0.249)		
Other race		0.409		0.086	(-0.231, 0.403)			0.081	(-0.395, 0.557)		
Other Hispanic		0.355		0.811	(0.463, 1.158)			0.811	(0.322, 1.300)		
Time effect		1.000		0.978	(0.951, 1.005)			0.978	(0.947, 1.008)		
Quadratic link function											
Environmental factors											
a-Tocopherol	1	0.560	1	0.558	(0.500, 0.617)	31.2	1	0.526	(0.288, 0.764)	0.614	(0.565, 0.664)
g-tocopherol	2	0.490	2	0.499	(0.453, 0.544)	24.9	3	0.333	(0.176, 0.489)	0.389	(0.324, 0.454)
Retinyl-palmitate	3	0.420	3	0.422	(0.363, 0.482)	17.8	4	0.279	(0.090, 0.467)	0.325	(0.242, 0.409)
Retinol	7	0.140	7	0.167	(0.108, 0.227)	2.8	8	-0.006	(-0.181, 0.169)	-0.007	(-0.078, 0.064)
3,3,4,4,5-pncb	8	0.070	8	0.073	(0.024, 0.122)	0.5	6	0.216	(0.027, 0.405)	0.252	(0.164, 0.341)
PCB194	6	-0.210	6	-0.209	(-0.269, -0.149)	4.4	2	0.378	(0.137, 0.619)	0.441	(0.352, 0.531)
2.3.4.6.7.8.hxcdf	5	-0.280	5	-0.268	(-0.327, -0.209)	7.2	7	-0.061	(-0.221, 0.100)	-0.071	(-0.141, -0.001)
trans.b.carotene	4	-0.350	4	-0.335	(-0.388, -0.283)	11.3	5	-0.273	(-0.44, -0.106)	-0.319	(-0.381, -0.257)
Covariates											
Intercept		0.000		0.877	(0.653, 1.100)			2.202	(1.478, 2.925)		
Age		0.005		0.007	(0.004, 0.009)			-0.023	(-0.035, -0.011)		
Sex (female)		-0.076		-0.061	(-0.158, 0.036)			-0.150	(-0.465, 0.165)		
Ethnicity											
Non-Hispanic white		Ref		Ref				Ref			
Non-Hispanic black		-0.138		-0.078	(-0.206, 0.050)			-0.004	(-0.395, 0.387)		
Mexican American		0.175		0.219	(0.081, 0.358)			0.323	(-0.070, 0.717)		
Other race		0.409		0.642	(0.372, 0.911)			0.095	(-0.763, 0.953)		
Other Hispanic		0.355		0.152	(-0.093, 0.397)			-0.125	(-0.976, 0.726)		
Time effect		1.000		1.014	(0.987, 1.041)			1.013	(0.983, 1.044)		

704

705

706

707

708

709

710

711 **Figure titles and legends**

712 **Fig. 1** Correlation matrix of Pearson correlation coefficients of 8 factors and triglycerides in
713 NHANES 2002-2003 (N=800).

714 **Fig. 2** Results from PLSI linear regression and multivariable linear regression in NHANES
715 2002-2003 (d=8, N=800). Bars show the estimated relative importance (absolute value of
716 estimated coefficient) of 8 environmental factors on continuous triglycerides. Red/green color
717 represents positive/negative effect. Error bars indicate 95% CIs.

718 **Fig. 3** Results from PLSI quantile regression and multivariable quantile regression in NHANES
719 2002-2003 (d=8, N=800). Bars show the estimated relative importance (absolute value of
720 estimated coefficient) of 8 environmental factors on three quartiles (25th, 50th, and 75th
721 percentiles) of triglycerides. Red/green color represents positive/negative effect. Error bars
722 indicate 95% CIs.

723 **Fig. 4** Results from PLSI logistic regression and multivariable logistic regression in NHANES
724 2002-2003 (N=800). Bars show the estimated relative importance (absolute value of estimated
725 coefficient) of 8 environmental factors on dichotomized triglycerides. Red/green color represents
726 positive/negative effect. Error bars indicate 95% CIs.

Figures

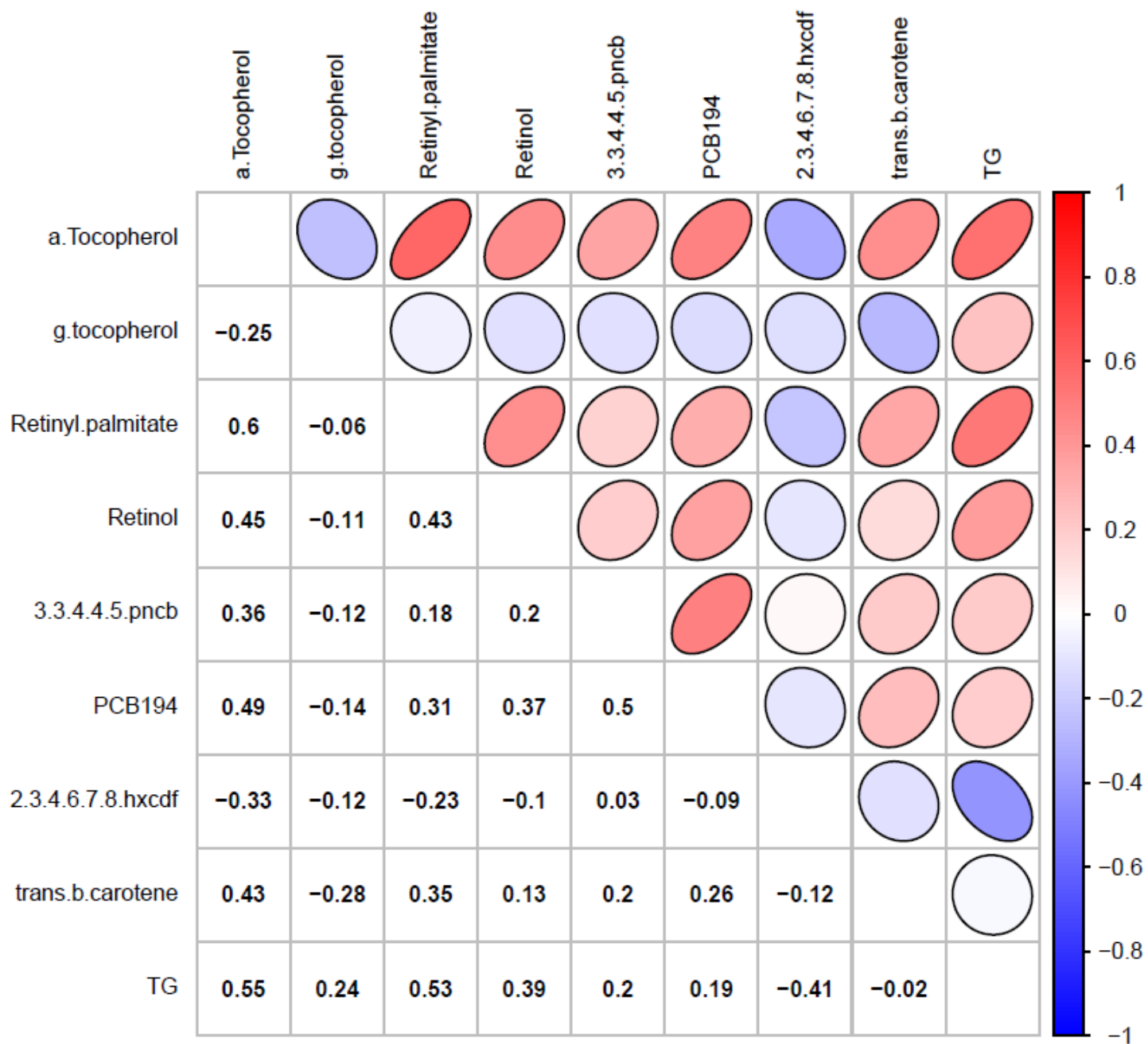


Figure 1

Correlation matrix of Pearson correlation coefficients of 8 factors and triglycerides in NHANES 2002-2003 (N=800).

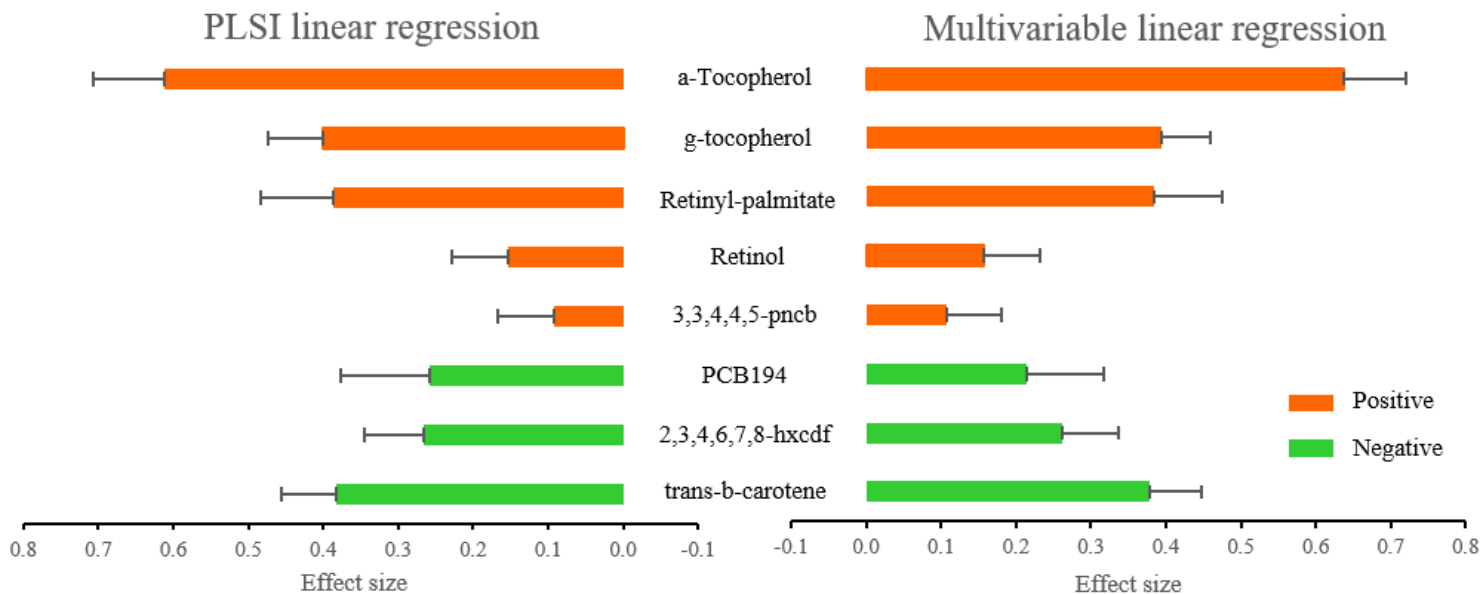


Figure 2

Results from PLSI linear regression and multivariable linear regression in NHANES 2002-2003 (d=8, N=800). Bars show the estimated relative importance (absolute value of estimated coefficient) of 8 environmental factors on continuous triglycerides. Red/green color represents positive/negative effect. Error bars indicate 95% CIs.

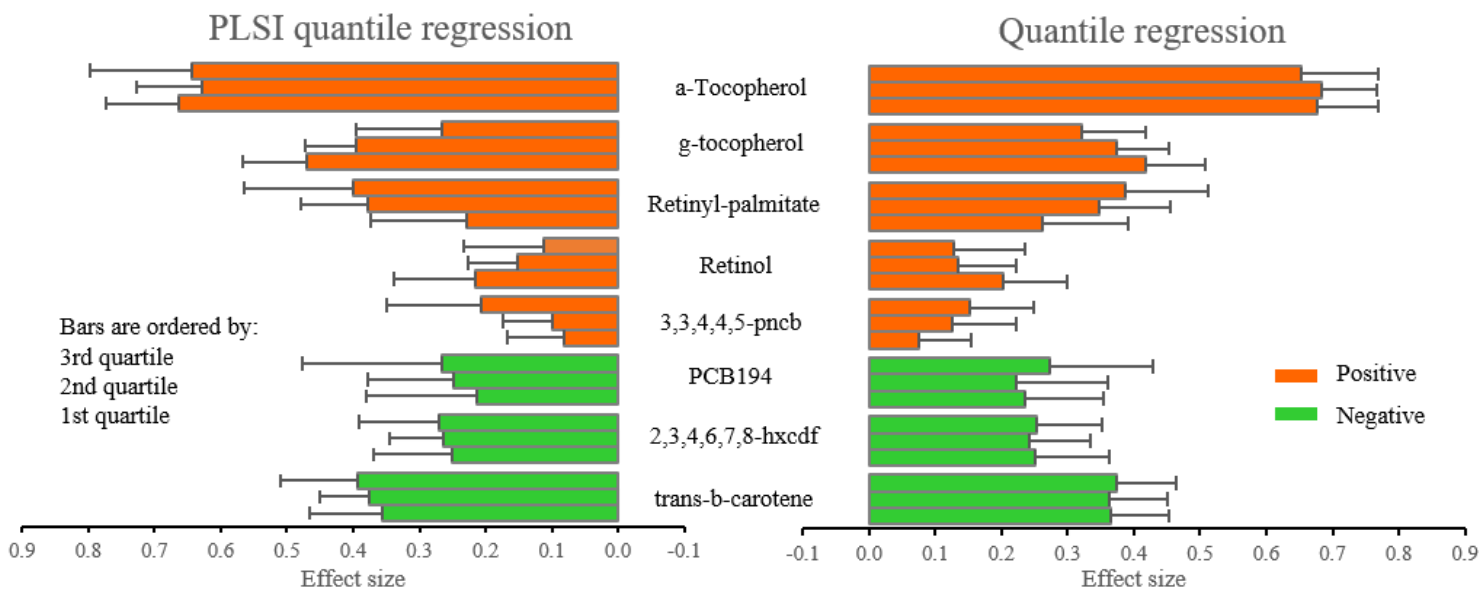


Figure 3

Results from PLSI quantile regression and multivariable quantile regression in NHANES 2002-2003 (d=8, N=800). Bars show the estimated relative importance (absolute value of estimated coefficient) of 8 environmental factors on three quartiles (25th, 50th, and 75th percentiles) of triglycerides. Red/green color represents positive/negative effect. Error bars indicate 95% CIs.

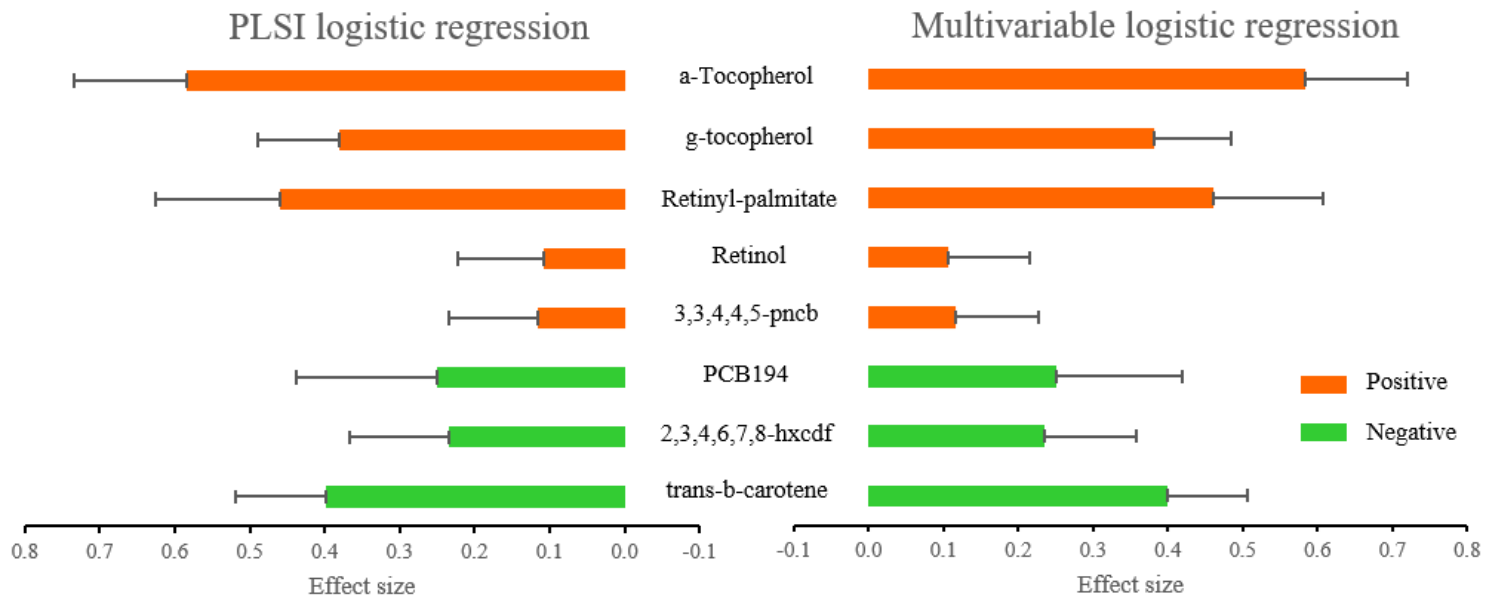


Figure 4

Results from PLSI logistic regression and multivariable logistic regression in NHANES 2002-2003 (N=800). Bars show the estimated relative importance (absolute value of estimated coefficient) of 8 environmental factors on dichotomized triglycerides. Red/green color represents positive/negative effect. Error bars indicate 95% CIs.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Additionalfile2.csv](#)
- [Additionalfile1tracked.docx](#)
- [Additionalfile3.Rmd](#)