

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 Wang<sup>1</sup>, Yinxiang Wu<sup>1</sup>, Melanie Jacobson<sup>3</sup>, Myeonggyun Lee<sup>1</sup>, Peng Jin<sup>1</sup>, Leonardo**  
7 **Trasande<sup>1,2,3</sup>, Mengling Liu<sup>1,2,\*</sup>**

8 <sup>1</sup> Department of Population Health, NYU Langone Health, New York, NY, USA

9 <sup>2</sup> Department of Environmental Medicine, NYU Langone Health, New York, NY, USA

10 <sup>3</sup> Department of Pediatrics, Divisions of Nephrology and Environmental Pediatrics, NYU  
11 Langone Health, New York, NY, USA

12 Yuyan Wang: [yuyan.wang@nyulangone.org](mailto:yuyan.wang@nyulangone.org);

13 Yinxiang Wu: [yinxiang.wu@nyulangone.org](mailto:yinxiang.wu@nyulangone.org);

14 Melanie Jacobson: [melanie.jacobson2@nyulangone.org](mailto:melanie.jacobson2@nyulangone.org);

15 Myeonggyun Lee: [myeonggyun.lee@nyulangone.org](mailto:myeonggyun.lee@nyulangone.org);

16 Peng Jin: [peng.jin@nyulangone.org](mailto:peng.jin@nyulangone.org);

17 Leonardo Trasande: [leonardo.trasande@nyulangone.org](mailto:leonardo.trasande@nyulangone.org);

18

19 **\*Corresponding author**

20 **Mengling Liu, PhD**

21 Email address: [mengling.liu@nyulangone.org](mailto: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 exposure factors 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 2012–2017 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. In each run of analysis, the WQS method assumes that all  
86 exposures are associated with the outcome in one direction, and then derives a one-dimensional  
87 weighted sum score of the exposures under the assumed direction for the estimation of overall  
88 effect. BKMR estimates the posterior inclusion probability (PIP) as the measure of importance  
89 for environmental exposures using a flexible nonparametric Bayesian framework. However, the  
90 estimated PIPs sometimes can be close to 1 for exposures that have strong effects on the  
91 outcome, and thus do not directly provide useful information on effect direction and relative  
92 importance. In addition, WQS and BKMR have been generalized to study environmental  
93 mixtures with several types of outcomes, such as WQS for longitudinal outcomes (13) and

94 BKMR for time-to-event outcomes (14). However, a general modeling framework that can  
 95 alleviate the above limitations in environmental health research is still desired (15).

96 Partial-linear single-index (PLSI) models are a family of semiparametric models that reside  
 97 between the completely unstructured nonparametric models and restrictive parametric regression  
 98 models (16-18). By reducing multiple exposures into the single index, the PLSI models can  
 99 reduce the “curse of dimensionality” issue and improve modeling efficiency. The application and  
 100 performance of single-index linear regression for analysis of environmental exposures with  
 101 continuous outcomes has been evaluated previously (pending publication). Specifically, the PLSI  
 102 modeling framework allows the associations between exposures and outcomes to be in the  
 103 positive or negative direction, provides explicit and interpretable quantification on the relative  
 104 direction and importance of the exposures, and models these effects with flexibility. In recent  
 105 years, research on PLSI models has attracted increasing attention and extended to different types  
 106 of outcomes, such as categorical (19-21), time-to-event (22-25) and longitudinal (26-29)  
 107 outcomes. Table 1 summarizes the outcome types of interest and corresponding PLSI models  
 108 with key references and their corresponding counterpart parametric models.

109 **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	(16), (19), (20), (30),(31), (32), (33), (34), (35), (36)	(1)
	PLSI quantile regression	Quantile regression	(37), (38), (39), (40), (41), (42)	(2)
Categorical (binary)	PLSI generalized linear (logistic) regression	Generalized linear (logistic) regression	(16), (20), (34), (36)	(3)
Time-to-event	PLSI PH model	Cox PH model	(22), (23), (24), (25)	(4)
Longitudinal	PLSI mixed-effects model	Linear mixed-effects model	(26), (27), (43), (44), (45)	(5)

110 The main goal of this study was to unify the resource advantages of PLSI models into one  
 111 general framework for analyzing environmental factors, and to demonstrate their values in  
 112 environmental research for different types of health outcomes. We exemplified the use of PLSI

113 models in assessing the associations between correlated environmental factors with health  
 114 outcomes using real and simulated datasets based on National Health and Nutrition Examination  
 115 Survey (NHANES) 2003-2004 cycle. Another aim was to develop effective computation  
 116 algorithms for the PLSI models and to consolidate these models using R packages.

## 117 **Methods**

### 118 NHANES dataset

119 To demonstrate the PLSI models, we used the data from the NHANES 2003-2004 cycle based on  
 120 the original paper by Patel et al (46), which systematically evaluated the associations of  
 121 environmental factors with serum lipid levels. We used serum triglyceride concentrations as the  
 122 primary outcome for demonstration and also considered three demographic variables, age, sex,  
 123 and race/ethnicity as potential confounders. Participants with data on serum triglycerides,  
 124 environmental factors and confounders were included in this study (n=800). Details on data pre-  
 125 processing are provided in [Additional file 1: Figure S1](#). Subjects provided written informed  
 126 consent, and the Institutional Review Board of the National Center for Health Statistics approved  
 127 the survey (47). [Table 2](#) summarizes the final variables included in analyses, and [Figure 1](#) shows  
 128 the correlation matrix of the final 8 environmental factors and triglycerides. The dataset is  
 129 provided as [Additional file 2](#), and the R codes conducting data cleaning is included in the R  
 130 markdown file ([Additional file 3](#)).

131 **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

Confounders	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
	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

132 For notational convention throughout this article, we let  $Y$  denote the outcome,  
 133  $X = (X_1, \dots, X_8)$  denote the 8 exposure variables to be modeled into the “single index” term, and  
 134 vector  $Z$  represent the confounders (age, sex, and race/ethnicity). The outcome, continuous  
 135 triglycerides, and all exposure variables, except for retinol, were log-transformed, and all  
 136 exposure variables were standardized to have mean of zero and standard deviations of 1 before  
 137 model fitting. We use  $\beta$ 's to denote the single index coefficients that characterize the relative  
 138 direction and importance of each exposure  $X_i$ , and  $\gamma$  for the corresponding linear coefficient  
 139 vector for confounder vector  $Z$ . To ensure model identifiability, the  $L_2$  norm of  $\beta$ 's  
 140 (i.e.  $\sqrt{\beta_1^2 + \dots + \beta_8^2}$ ) is set to be 1 with the first component  $\beta_1 > 0$ , which are the standard  
 141 parametrization constraints for all PLSI models.

142 Continuous outcome: mean regression

143 The PLSI linear regression model is considered as a generalization of both standard linear  
 144 regression and missing-link function problem in linear modeling (48), and specified as

$$145 \quad Y = g\left(\sum_{j=1}^8 \beta_j X_j\right) + \gamma'Z + \varepsilon \quad (1)$$

146 The semiparametric PLSI linear regression has the parametric component  $\sum_{j=1}^8 \beta_j X_j$  and  $\gamma'Z$  for  
 147 easy linear representation and interpretation, and the nonparametric components  $g(\cdot)$  is totally

148 unspecified and represents the overall effect of single index. When the estimated  $g(\cdot)$  is  
149 monotone, the effect of  $X_j$  can be interpreted qualitatively using the sign of  $\beta_j$ . If  $g(\cdot)$  is  
150 monotone increasing, then a positive sign for  $\beta_j$  suggests increased conditional expectation of  $Y$   
151 at larger value of  $X_j$ , and vice versa for a negative sign. As the overall scale of  $\beta$  is set,  $|\beta_j|$  can  
152 be explained as the relative importance of  $X_j$  affecting the mean of outcome  $Y$  as  $X_j$  is perturbed  
153 while  $g(\cdot)$  and other variables are held fixed. We can also intuitively interpret  $\beta_j^2$  as the  
154 proportion of contribution to the single index by variable  $X_j$  because, when  $(X_1, X_2, \dots, X_8)$  are  
155 independent,  $\beta_j^2$  simply represents  $X_j$ 's variance contribution.

#### 156 Continuous outcome: quantile regression

157 Beyond the commonly-considered effects of environmental factors on the mean of a continuous  
158 outcome, sometimes we are interested in the specific relations cross multiple points of the  
159 outcome's distribution, such as higher quantiles of triglycerides (49), higher quantiles of blood  
160 pressure (50), low quantiles of birth weight (51), or lower quantiles of intelligence quotient  
161 scores (52). Moreover, when the distribution of continuous outcome deviates from Gaussian,  
162 modeling the median can be more robust than evaluating the mean by conventional linear  
163 regression (53). For this purpose, quantile regression (QR), which was originally proposed by  
164 Koenker and Bassett (54) and used as a useful technique in econometrics (55) and growth curve  
165 analysis (56), enables us to study the associations of environmental factors with continuous  
166 health outcomes as various quantiles across its distribution. PLSI quantile regression is a  
167 combination of the PLSI technique and QR (40, 41), and thus we consider it for the analysis of  
168 joint effects of multiple environmental factors on the quantile(s) of continuous outcome variable.

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

171 
$$\theta_{\tau}(Y|X, Z) = g_{\tau}\left(\sum_{j=1}^8 \beta_{\tau j} X_j\right) + \gamma'_{\tau} Z \quad (2)$$

172 Interpretation of coefficients  $\beta_{\tau}$ 's in the PLSI quantile regression is similar to that of PLSI linear  
 173 regression, with the difference being that the associations are now with the conditional quantiles  
 174 of outcome variable  $\theta_{\tau}(Y|X, Z)$  instead of the mean.

175 **Categorical outcome: generalized linear regression**

176 PLSI generalized linear regression can be employed for categorical outcomes, such as binary,  
 177 multinomial, or count variables. Here we considered the binary outcome of high triglycerides (>  
 178 150 milligrams per deciliter) (57), which accounted for 30.75% of the 800 subjects. The PLSI  
 179 logistic model is specified as

180 
$$\text{logit}(P(Y = 1|X, Z)) = g\left(\sum_{j=1}^8 \beta_j X_j\right) + \gamma' Z \quad (3)$$

181 The interpretation of coefficients is based on the log odds that response value is '1' conditioning  
 182 on the predictors, and  $\beta_j$  represents the relative direction and importance of  $X_j$  associated with  
 183 the log odds of high triglycerides when scale of  $\beta$  is set and  $g(\cdot)$  and other variables are held  
 184 fixed. The logit function can be adapted accordingly to the type of categorical outcome.

185 **Time-to-event outcome: proportional hazards model**

186 The Cox proportional hazards (PH) regression has been the pivotal model in time-to-event  
 187 analysis since Sir Cox proposed it in 1972 (58, 59). The Cox PH regression models the hazard  
 188 function and assumes that covariates have linear effects on the log hazard function. Combining  
 189 PLSI modeling technique and Cox PH regression, the PLSI PH model is specified as

190 
$$\lambda(t|X, Z) = \lambda_0(t) \exp \left\{ g \left( \sum_{j=1}^8 \beta_j X_j \right) + \gamma' Z \right\}, \quad (4)$$

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

193 **Longitudinal outcome: mixed-effects model**

194 Longitudinal studies arise frequently in environmental research, in which outcomes are measured  
 195 repeatedly over a period of time with either baseline or time-dependent environmental factors.

196 As measurements from the same subject are often correlated, subject-specific random effects are  
 197 used to accommodate within-subject dependence and to explain across-subject heterogeneity.

198 Mixed-effects models provide a general and flexible framework for modeling longitudinal data,

199 consisting of two modeling components: fixed effects and random effects, characterizing the

200 population mean and individual variation, respectively (60, 61). Mixed-effects models in general

201 are amenable to missing data and can accommodate missing completely at random or missing at

202 random (60, 62). Without loss of generality, we consider a longitudinal study with  $N$  subjects

203 and the  $i$ th subject has  $n_i$  observations over time. Repeated measures of the outcome are denoted

204 by  $Y_{ij}$ , exposure vector  $X_{ij}$ , covariate vector  $Z_{ij}$  and observation time  $T_{ij}$ , and then the observed

205 full dataset is  $\{(Y_{ij}, X_{ij}, Z_{ij}, T_{ij}), i = 1, \dots, N, j = 1, \dots, n_i\}$ .

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

207 
$$Y_{ij} = g \left( \sum_{l=1}^8 \beta_l X_{ijl} \right) + Z'_{ij} \gamma + b_i + \omega T_{ij} + \varepsilon_{ij}, \quad (5)$$

208 where  $b_i$  represents the subject-specific random intercept and  $\omega$  represents the time effect on the

209 outcome. Note that PLSI mixed-effects model can accommodate additional random effects and

210 other model specifications of fixed effects and interactions. The index coefficient  $\beta_l$  can be

211 explained as the relative direction and importance of  $X_{ijl}$  as  $X_{ijl}$  is perturbed when scale of  $\beta$  is  
 212 set and  $g(\cdot)$  and other variables are held fixed, and  $g(\cdot)$  represents the overall effect of the single  
 213 index with the mean of longitudinal outcome.

#### 214 Simulation settings

215 Since the NHANES survey dataset does not have time-to-event outcome nor longitudinal  
 216 outcome, we conducted simulations to demonstrate the PLSI PH model and PLSI mixed-effects  
 217 model. The coefficients for the 8 environmental factors and three confounding variables were set  
 218 based on the results from the PLSI linear regression for continuous triglycerides. We kept the  
 219 original direction of these associations and the absolute rank for each environment factor, and set  
 220 the effect sizes in a wider range to be more distinguishable (see details in [Table 3](#) and [Table 4](#)).

221 Moreover, we considered the link function  $g(\cdot)$  to be either  $g(x) = x$  to facilitate the direct  
 222 comparison with the parametric models, or as a quadratic function  $g(x) = x^2$  to mimic the  
 223 scenario with nonlinear effects and pair-wise interactions between the exposures as

$$224 \quad g\left(\sum_{j=1}^8 \beta_j X_j\right) = \beta_1^2 X_1^2 + \dots + \beta_8^2 X_8^2 + 2\beta_1\beta_2 X_1 X_2 + \dots + 2\beta_7\beta_8 X_7 X_8.$$

225 Time-to-event outcomes were generated using model (4) with  $\lambda_0 = 1$  in identity link  
 226 function scenario and  $\lambda_0 = 1/\exp(2)$  in quadratic link function scenario, and censoring rate as  
 227 20%. Longitudinal outcomes were generated using model (5) with  $t_{ij}$  ranged  $[1, 6]$  and  $\omega = 1$ .

228 The number of possible observations for each subject was assumed to vary randomly between 2  
 229 and 6. The errors followed a first order autoregressive process (i.e. AR(1)), with the  
 230 autocorrelation as 0.4 and standard deviation as 1.5 to mimic decreasing dependence with time.

231 All details of data generation used in these simulations are included in the R markdown file  
 232 ([Additional file 3](#)).

#### 233 Performance evaluation

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

#### 243 **Statistical software**

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

#### 251 **Results**

##### 252 **Continuous triglycerides: PLSI mean regression**

253 We applied the PLSI linear regression and multivariable linear regression to study the  
254 associations of the 8 environmental factors with continuous triglycerides, and summarized the  
255 estimates in [Figure 2](#) (numerical results in [Additional file 1: Table S1](#)). The ranks, estimated  
256 coefficients, and directions were similar between these two models, and the estimated link

257 function was close to be linear ([Additional file 1: Figure S3](#)). As the estimated link function was  
258 monotone and increasing, the positive estimates indicated a positive association with  
259 triglycerides. Specifically, a-Tocopherol had a  $\hat{\beta}_1 = 0.612$  and 95% CI of (0.517, 0.707),  
260 indicating that a-Tocopherol had the strongest positive association with triglycerides among the 8  
261 factors, and made about 37.4% contribution to the single index; trans-b-carotene had the most  
262 negative association of  $\hat{\beta}_8 = -0.383$ . These results were consistent with original results from  
263 Patel's study, which also observed a-Tocopherol with the strongest positive and trans-b-carotene  
264 with the strongest negative association with triglycerides (46). As the 8 environmental factors  
265 showed both positive and negative associations with triglycerides, this application highlighted  
266 the need of statistical methods to accommodate both directional effects for studying multiple  
267 environmental exposures.

#### 268 Continuous triglycerides: PLSI quantile regression

269 We applied the PLSI quantile regression to study the associations between 8 exposures and  
270 three quartiles (25<sup>th</sup>, 50<sup>th</sup>, and 75<sup>th</sup> percentiles) of triglycerides and summarized the main results  
271 in [Figure 3](#) (numerical results in [Additional file 1: Table S2](#)). We observed that the estimated  
272 link functions for all three quartiles were increasing and close to linear ([Additional file 1: Figure](#)  
273 [S4](#)), which explained the similarities between the results of the PLSI quantile regressions and  
274 regular quantile regressions. In addition, the 8 environmental factors showed fairly consistent  
275 associations across the three quartiles of triglycerides. For example, a-Tocopherol was the factor  
276 having the strongest positive association with triglycerides and trans-b-carotene was the factor  
277 having the strongest negative association with triglycerides at all three quartiles.

#### 278 Binary triglycerides: PLSI logistic regression

279 For dichotomized triglycerides, the ranks and estimates from PLSI logistic regression and  
280 multivariable logistic regression are shown in [Figure 4](#) (numerical results in [Additional file 1:](#)  
281 [Table S3](#)), which demonstrated similar results from these two models. The estimated link  
282 function by PLSI logistic regression was monotone increasing and close to be linear ([Additional](#)  
283 [file 1: Figure S5](#)). Thus, the estimated directions can be interpreted qualitatively and the  
284 estimated coefficients represented the relative importance of each exposure on the log odds of  
285 high triglycerides. For example, the estimated coefficient of a-Tocopherol was  $\hat{\beta}_1 = 0.584$  (95%  
286 CI: 0.433-0.735), which represented that a-Tocopherol had the strongest positive association  
287 with the odds of high triglycerides among the 8 factors.

#### 288 Simulated time-to-event outcome: PLSI PH model

289 We summarize the simulation results from both PLSI PH model and Cox PH model in [Table](#)  
290 [3](#). Under the identity link function setting, results from the PLSI PH model and the conventional  
291 Cox PH model were very similar as expected, and both close to the true values. The PLSI PH  
292 model estimated the link function to be close to the true linear function ([Additional file 1: Figure](#)  
293 [S6 \(a\)](#)). Under the quadratic link function setting, results from the PLSI PH model were still  
294 consistent to true coefficients, but the conventional Cox PH model failed for most of the  
295 environmental factors because the linear model assumption was insufficient. The PLSI PH model  
296 also captured the U-shape and estimated the link function close to the true quadratic function  
297 ([Additional file 1: Figure S6 \(b\)](#)).

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

#### 299 Simulated longitudinal outcome: PLSI mixed-effects model

300 The results from PLSI mixed-effects model and linear mixed-effects model under identify or  
301 quadratic link function are presented in [Table 4](#). Under the identity link function setting, the

302 PLSI mixed-effects model estimated all coefficients close to the true coefficients with correct  
303 directions, and conventional linear mixed-effects model also had similar estimations. The  
304 estimated link function by PLSI mixed-effects model was close to the true linear function  
305 (Additional file 1: Figure S7 (a)). Under the quadratic link function setting, the results from PLSI  
306 mixed-effects model were still consistent; however, the conventional linear mixed-effects model  
307 clearly showed biased results for some factors like PCB194. The estimated link function by PLSI  
308 mixed-effects model had a U-shape and was close to the true quadratic function (Additional file  
309 1: Figure S7 (b)).  
310 (Table 4 should appear here)

## 311 Discussion

312 We presented five PLSI models aiming to provide a unified family of statistical models to assess  
313 the joint effects of environmental exposures on four types of health outcomes: continuous,  
314 categorical, time-to-event, and longitudinal outcomes. We demonstrated the flexibility and  
315 effectiveness of this PLSI family for modeling various types of outcomes using NHANES data  
316 supplemented with simulations. One contribution of this work is that the novel modeling options  
317 under the PLSI framework complement existing methods and address some common statistical  
318 challenges in the analysis of multiple environmental exposures, such as mixed directions,  
319 interactions, and non-linear effects. Another contribution is that coherent computation algorithms  
320 are developed for all the PLSI models and implemented using the existing R packages, which  
321 can facilitate direct applications in practice and reproducible research.

322 In our analyses of the cross-sectional NHANES studies for continuous and binary  
323 triglycerides by PLSI models, we found that the 8 environmental factors exhibited mixed  
324 directional associations with the outcome, with a-Tocopherol having the strongest positive

325 association and trans-b-carotene having the strongest negative association with triglycerides. A-  
326 Tocopherol and carotenes are transported in serum with HDL and LDL, and the level of serum a-  
327 Tocopherol depends on serum lipids (63, 64). The strong positive association between a-  
328 Tocopherol and triglycerides is expected (46), and the negative association between b-carotene  
329 and triglycerides is supported by previous studies (65, 66). Our results were consistent with the  
330 results of previously known and validated environmental chemical factors correlated with  
331 triglycerides (46), clearly demonstrating the value of PLSI models as a flexible and useful tool  
332 for analyzing complex exposures. Using additional simulations for time-to-event and  
333 longitudinal outcomes, we showed that the PLSI models could correctly identify the directions  
334 and magnitudes of associations for these environmental factors in scenarios with different types  
335 of outcomes.

336 In our NHANES applications of studying triglycerides continuously and categorically, we  
337 estimated that the link functions of PLSI models were very close to be linear, which were also  
338 reflected by the similar results with their counterpart parametric models. In general, standard  
339 errors from the PLSI models were larger than those from their counterpart parametric models,  
340 which was expected as the former are semiparametric models.

341 Moreover, our results remained consistent when we conducted several sensitivity analyses.  
342 First, we conducted a sensitivity analysis including all 22 environmental factors to investigate the  
343 performance of PLSI linear regression to handle highly correlated exposures. Results showed  
344 that the key observations on the important environmental factors were similar ([Additional file 1:  
345 Table S4](#)). When there are many highly correlated exposure factors ( $r > 0.9$ , [Additional file 1:  
346 Figure S2](#)), we also recommend to use p-values to rank the importance of variables because the  
347 value of estimates can be affected by collinearity. As shown in [Table S4](#), the 8 selected  
348 environmental factors still showed top ranks among the 22 factors, except for PCB194 which

349 was highly correlated with other PCBs. In addition, we conducted another sensitivity weighted  
350 analysis incorporating the laboratory subsample C weights from NHANES 2003-2004 cycle  
351 (following general guideline to use the weights from “least common denominator”) (67), and the  
352 weighted results ([Additional file 1: Table S5](#)) were similar with the results from unweighted  
353 models. Note that most of the PLSI models are readily incorporate weights in R function codes  
354 ([Additional file 3](#)).

355 Interaction among multiple correlated environmental factors is very common, and it has  
356 been long appreciated that the co-exposures may have synergistic (additive or multiplicative) or  
357 antagonistic effects on health outcomes (69). For parametric models, it’s difficult to directly  
358 model the interaction effects among co-exposures if we don’t know the ‘degree of interaction’.  
359 However, PLSI models can handle the interaction easily through the unknown link function as  
360 we evaluated using the simulations. Specifically, in our simulated time-to-event and longitudinal  
361 analyses with quadratic link function, which reflected both the pairwise interactions and non-  
362 linear quadratic effects, both PLSI PH model and PLSI mixed-effects were able to capture the U-  
363 shape link function and correct direction and importance of the environmental factors, while  
364 parametric models failed in most factors because the parametric assumptions were no longer  
365 satisfied. Therefore, PLSI models readily accommodate the factors showing non-linear or  
366 interactive effect on the health outcome.

367 There are other ways and models using various definitions of weighted sums to model the  
368 joint effect for multiple environmental components. For example, molar sums were used to show  
369 relationships between prenatal phenol and phthalate exposures and birth outcome (70), and a  
370 potency-weighted sum was used to calculate phthalates exposures among reproductive-aged  
371 women (71). The weights for environmental factors can be calculated from their expected  
372 potency relative to a reference factor, like the common cases in toxicology (72), or based on their

373 percent contribution to the total mixture effect, like WQS (9). PLSI models can be considered as  
374 one of these weighting approaches, and their advantages from the semiparametric structure are  
375 evident compared with existing methods, especially for the scenarios when the environmental  
376 exposures have mixed-directional associations and/or a potential high-degree interaction.  
377 Meanwhile, due to the flexibility of the nonparametric link function, PLSI models can represent  
378 complex joint effects more than additive structures (73), which is commonly encountered since  
379 environmental exposures may act together in a biological sense via a shared mechanistic  
380 pathway (4). The ability of handling various types of outcomes is another important advantage of  
381 the proposed PLSI framework. This is important because, with the accumulation of  
382 environmental exposure measurements and development of data collection methods, time-to-  
383 event or longitudinal studies are desired to explore the associations over time.

384 In this study, the coherent algorithms for PLSI models are based on the ‘gam’ and ‘gamm’  
385 functions from ‘mgcv’ package and ‘qgam’ function from ‘qgam’ package in R, which includes  
386 many of the generalized additive model (GAM) fitting techniques developed by Simon Wood et  
387 al (74). The rationale behind this algorithms is to use ‘gam’, ‘qgam’ or ‘gamm’ call (usually  
388 using penalized regression splines or similar smoothers) to profile out the smooth model  
389 coefficients and smoothing parameters for estimation of the link function contained in PLSI  
390 model, leaving only a finite parameter vector to be estimated by a general purpose optimizer.  
391 Based on this algorithm, it is easy to adapt the models to include multiple single index terms,  
392 parametric terms, and further smoothing. We have compared the estimates for single index  
393 models among different iterative procedures using existing packages (e.g., projection pursuit  
394 regression with one term using ‘ppr’ function; ‘sim.est’ function from ‘simest’ package) in  
395 various simulations, and they have similar estimation performance. We finally chose ‘gam’ call  
396 series because of its flexibility for covariate adjustment and ability of modeling various types of

397 outcomes. This ‘gam’, ‘qgam’, ‘gamm’ call approach has demonstrated efficient and robust  
398 performance in our numerical studies, and we believe this coherent algorithm strategy wrapped  
399 as a toolbox is beneficial for practical application.

400 The PLSI models considered here may not be directly applicable to extreme high-  
401 dimensional settings, for which we could consider using extensions with adaptive LASSO (75),  
402 smoothly clipped absolute deviation penalty (76), and smooth-threshold estimating equations  
403 (77). Another future research direction is to extend from the single index to multiple-index  
404 models, such as the projection pursuit regression (78), so that more complex data structures and  
405 exposure effect patterns can be captured and modeled.

## 406 **Conclusions**

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

## 410 **Additional files**

411 **Addition file 1: Figure S1.** Data flow diagram for deriving 800 subjects and 8 environmental  
412 factors. **Figure S2.** Correlation matrix of Pearson correlation coefficient of 22 factors and  
413 triglycerides in NHANES 2002-2003 (N=800). **Table S1.** Results from PLSI linear regression and  
414 multivariable linear regression in NHANES 2002-2003. **Figure S3.** Estimated link function by PLSI  
415 linear regression in NHANES 2002-2003. **Tables S2.1-S2.3.** Results from PLSI quantile regressions  
416 and multivariable quantile regression at three quantiles (25th, 50th, and 75th percentiles) of  
417 triglycerides in NHANES 2002-2003. **Figure S4.** Estimated link functions by PLSI quantile  
418 regressions at three quartiles in NHANES 2002-2003. (a) 25th percentile; (b) 50th percentile; (c)

419 75th percentile. **Table S3.** Results from PLSI logistic regression and multivariable logistic  
 420 regression in NHANES 2002. **Figure S5.** Estimated link function by PLSI logistic regression in  
 421 NHANES 2002-2003. **Figure S6.** Estimated link functions by PLSI PH model in simulated time-to-  
 422 event study. (a) identity link function; (b) quadratic link function. **Figure S7.** Estimated link  
 423 functions by PLSI mixed-effects model in simulated longitudinal study. (a) identity link function;  
 424 (b) quadratic link function. **Table S4.** Sensitivity analysis results from PLSI linear regression and  
 425 multivariable linear regression in NHANES 2002-2003 with 22 environmental factors. **Table S5.**  
 426 Sensitivity analysis results from weighted PLSI linear regression and weighted linear regression in  
 427 NHANES 2002-2003 using NHANES laboratory subsample C weights.

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

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

### 433 **Abbreviations**

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

### 438 **Declarations**

439 **Ethics approval and consent to participate**

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

#### 442 **Consent for publications**

443 Not applicable.

#### 444 **Availability of data and materials**

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

#### 447 **Competing interests**

448 The authors declare that they have no competing interests.

#### 449 **Funding**

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

#### 452 **Authors' contributions**

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

#### 458 **Acknowledgements**

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

460

## 461 References

- 462 1. Wild CP. Complementing the genome with an "exposome": The outstanding challenge of  
 463 environmental exposure measurement in molecular epidemiology. *Cancer Epidemiology*  
 464 *Biomarkers & Prevention*. 2005;14(8):1847-50.
- 465 2. Stafoggia M, Breitner S, Hampel R, Basagana X. Statistical Approaches to Address Multi-  
 466 Pollutant Mixtures and Multiple Exposures: the State of the Science. *Curr Environ Health*  
 467 *Rep*. 2017;4(4):481-90.
- 468 3. Sanders AP, Claus Henn B, Wright RO. Perinatal and Childhood Exposure to Cadmium,  
 469 Manganese, and Metal Mixtures and Effects on Cognition and Behavior: A Review of Recent  
 470 Literature. *Curr Environ Health Rep*. 2015;2(3):284-94.
- 471 4. Hamra GB, Buckley JP. Environmental exposure mixtures: questions and methods to address  
 472 them. *Curr Epidemiol Rep*. 2018;5(2):160-5.
- 473 5. Carlin DJ, Rider CV, Woychik R, Birnbaum LS. Unraveling the health effects of  
 474 environmental mixtures: an NIEHS priority. *Environ Health Perspect*. 2013;121(1):A6-8.
- 475 6. Billionnet C, Sherrill D, Annesi-Maesano I, Study G. Estimating the Health Effects of  
 476 Exposure to Multi-Pollutant Mixture. *Annals of Epidemiology*. 2012;22(2):126-41.
- 477 7. Mann RM, Hyne RV, Choung CB, Wilson SP. Amphibians and agricultural chemicals:  
 478 review of the risks in a complex environment. *Environ Pollut*. 2009;157(11):2903-27.
- 479 8. Chaumont A, Nickmilder M, Dumont X, Lundh T, Skerfving S, Bernard A. Associations  
 480 between proteins and heavy metals in urine at low environmental exposures: Evidence of  
 481 reverse causality. *Toxicology Letters*. 2012;210(3):345-52.
- 482 9. Carrico C, Gennings C, Wheeler DC, Factor-Litvak P. Characterization of Weighted  
 483 Quantile Sum Regression for Highly Correlated Data in a Risk Analysis Setting. *Journal of*  
 484 *Agricultural Biological and Environmental Statistics*. 2015;20(1):100-20.

- 485 10. Czarnota J, Gennings C, Colt JS, De Roos AJ, Cerhan JR, Severson RK, et al. Analysis of  
486 Environmental Chemical Mixtures and Non-Hodgkin Lymphoma Risk in the NCI-SEER  
487 NHL Study. *Environmental Health Perspectives*. 2015;123(10):965-70.
- 488 11. Bobb JF, Valeri L, Claus Henn B, Christiani DC, Wright RO, Mazumdar M, et al. Bayesian  
489 kernel machine regression for estimating the health effects of multi-pollutant mixtures.  
490 *Biostatistics*. 2015;16(3):493-508.
- 491 12. Valeri L, Mazumdar MM, Bobb JF, Henn BC, Rodrigues E, Sharif OIA, et al. The Joint  
492 Effect of Prenatal Exposure to Metal Mixtures on Neurodevelopmental Outcomes at 20-40  
493 Months of Age: Evidence from Rural Bangladesh. *Environmental Health Perspectives*.  
494 2017;125(6).
- 495 13. Levin-Schwartz Y, Gennings C, Schnaas L, Del Carmen Hernandez Chavez M, Bellinger  
496 DC, Tellez-Rojo MM, et al. Time-varying associations between prenatal metal mixtures and  
497 rapid visual processing in children. *Environ Health*. 2019;18(1):92.
- 498 14. Zhang L, Kim I. Semiparametric Bayesian kernel survival model for evaluating pathway  
499 effects. *Statistical Methods in Medical Research*. 2019;28(10-11):3301-17.
- 500 15. Gibson EA, Nunez Y, Abuawad A, Zota AR, Renzetti S, Devick KL, et al. An overview of  
501 methods to address distinct research questions on environmental mixtures: an application to  
502 persistent organic pollutants and leukocyte telomere length. *Environmental Health*.  
503 2019;18(1).
- 504 16. Ichimura H. Semiparametric Least-Squares (Sls) and Weighted Sls Estimation of Single-  
505 Index Models. *Journal of Econometrics*. 1993;58(1-2):71-120.
- 506 17. Horowitz JL, Hardle W. Direct semiparametric estimation of single-index models with  
507 discrete covariates. *Journal of the American Statistical Association*. 1996;91(436):1632-40.

- 508 18. Wang JL, Xue LG, Zhu LX, Chong YS. Estimation for a Partial-Linear Single-Index Model.  
 509 Annals of Statistics. 2010;38(1):246-74.
- 510 19. Hardle W, Hall P, Ichimura H. Optimal Smoothing in Single-Index Models. Annals of  
 511 Statistics. 1993;21(1):157-78.
- 512 20. Carroll RJ, Fan JQ, Gijbels I, Wand MP. Generalized partially linear single-index models.  
 513 Journal of the American Statistical Association. 1997;92(438):477-89.
- 514 21. Yi GY, He WQ, Liang H. Analysis of correlated binary data under partially linear single-  
 515 index logistic models. Journal of Multivariate Analysis. 2009;100(2):278-90.
- 516 22. Wang W. Proportional hazards regression models with unknown link function and time-  
 517 dependent covariates. Statistica Sinica. 2004;14(3):885-905.
- 518 23. Huang JHZ, Liu LX. Polynomial spline estimation and inference of proportional hazards  
 519 regression models with flexible relative risk form. Biometrics. 2006;62(3):793-802.
- 520 24. Sun J, Kopciuk KA, Lu XW. Polynomial spline estimation of partially linear single-index  
 521 proportional hazards regression models. Computational Statistics & Data Analysis.  
 522 2008;53(1):176-88.
- 523 25. Li JB, Zhang RQ. Partially varying coefficient single index proportional hazards regression  
 524 models. Computational Statistics & Data Analysis. 2011;55(1):389-400.
- 525 26. Bai Y, Fung WK, Zhu ZY. Penalized quadratic inference functions for single-index models  
 526 with longitudinal data. J Multivariate Anal. 2009;100(1):152-61.
- 527 27. Li GR, Zhu LX, Xue LG, Feng SY. Empirical likelihood inference in partially linear single-  
 528 index models for longitudinal data. J Multivariate Anal. 2010;101(3):718-32.
- 529 28. Xu PR, Zhu LX. Estimation for a marginal generalized single-index longitudinal model.  
 530 Journal of Multivariate Analysis. 2012;105(1):285-99.

- 531 29. Zhao WH, Lian H, Liang H. GEE analysis for longitudinal single-index quantile regression. J  
 532 Stat Plan Infer. 2017;187:78-102.
- 533 30. Stoker TM. Consistent Estimation of Scaled Coefficients. *Econometrica*. 1986;54(6):1461-  
 534 81.
- 535 31. Hardle W, Stoker TM. Investigating Smooth Multiple-Regression by the Method of Average  
 536 Derivatives. *Journal of the American Statistical Association*. 1989;84(408):986-95.
- 537 32. Hardle W, Tsybakov AB. How Sensitive Are Average Derivatives. *Journal of Econometrics*.  
 538 1993;58(1-2):31-48.
- 539 33. Hristache M, Juditsky A, Spokoiny V. Direct estimation of the index coefficient in a single-  
 540 index model. *Annals of Statistics*. 2001;29(3):595-623.
- 541 34. Yu Y, Ruppert D. Penalized spline estimation for partially linear single-index models.  
 542 *Journal of the American Statistical Association*. 2002;97(460):1042-54.
- 543 35. Xia YC, Hardle W. Semi-parametric estimation of partially linear single-index models.  
 544 *Journal of Multivariate Analysis*. 2006;97(5):1162-84.
- 545 36. Liang H, Liu X, Li RZ, Tsai CL. Estimation and Testing for Partially Linear Single-Index  
 546 Models. *Annals of Statistics*. 2010;38(6):3811-36.
- 547 37. Chaudhuri P. Global Nonparametric-Estimation of Conditional Quantile Functions and Their  
 548 Derivatives. *Journal of Multivariate Analysis*. 1991;39(2):246-69.
- 549 38. Chaudhuri P, Doksum K, Samarov A. On average derivative quantile regression. *Annals of*  
 550 *Statistics*. 1997;25(2):715-44.
- 551 39. Wu TZ, Yu KM, Yu Y. Single-index quantile regression. *Journal of Multivariate Analysis*.  
 552 2010;101(7):1607-21.
- 553 40. Kong EF, Xia YC. A Single-Index Quantile Regression Model and Its Estimation.  
 554 *Econometric Theory*. 2012;28(4):730-68.

- 555 41. Lv YZ, Zhang RQ, Zhao WH, Liu JC. Quantile regression and variable selection of partial  
 556 linear single-index model. *Annals of the Institute of Statistical Mathematics*. 2015;67(2):375-  
 557 409.
- 558 42. Ma SJ, He XM. Inference for Single-Index Quantile Regression Models with Profile  
 559 Optimization. *Annals of Statistics*. 2016;44(3):1234-68.
- 560 43. Lai P, Li GR, Lian H. Quadratic inference functions for partially linear single-index models  
 561 with longitudinal data. *Journal of Multivariate Analysis*. 2013;118:115-27.
- 562 44. Li GR, Lai P, Lian H. Variable selection and estimation for partially linear single-index  
 563 models with longitudinal data. *Statistics and Computing*. 2015;25(3):579-93.
- 564 45. Li JB, Lian H, Jiang XJ, Song XY. Estimation and testing for time-varying quantile single-  
 565 index models with longitudinal data. *Computational Statistics & Data Analysis*. 2018;118:66-  
 566 83.
- 567 46. Patel CJ, Cullen MR, Ioannidis JPA, Butte AJ. Systematic evaluation of environmental  
 568 factors: persistent pollutants and nutrients correlated with serum lipid levels. *International*  
 569 *Journal of Epidemiology*. 2012;41(3):828-43.
- 570 47. Zipf G, Chiappa M, Porter KS, Ostchega Y, Lewis BG, Dostal J. National health and  
 571 nutrition examination survey: plan and operations, 1999-2010. *Vital Health Stat 1*.  
 572 2013(56):1-37.
- 573 48. Weisberg S, Welsh AH. Adapting for the Missing Link. *Annals of Statistics*.  
 574 1994;22(4):1674-700.
- 575 49. Di Angelantonio E, Sarwar N, Perry P, Kaptoge S, Ray KK, Thompson A, et al. Major  
 576 Lipids, Apolipoproteins, and Risk of Vascular Disease. *Jama-Journal of the American*  
 577 *Medical Association*. 2009;302(18):1993-2000.

- 578 50. Bind MA, Peters A, Koutrakis P, Coull B, Vokonas P, Schwartz J. Quantile Regression  
 579 Analysis of the Distributional Effects of Air Pollution on Blood Pressure, Heart Rate  
 580 Variability, Blood Lipids, and Biomarkers of Inflammation in Elderly American Men: The  
 581 Normative Aging Study. *Environmental Health Perspectives*. 2016;124(8):1189-98.
- 582 51. Burgette LF, Reiter JP, Miranda ML. Exploratory Quantile Regression With Many  
 583 Covariates An Application to Adverse Birth Outcomes. *Epidemiology*. 2011;22(6):859-66.
- 584 52. Ratcliff R, Thapar A, McKoon G. Individual differences, aging, and IQ in two-choice tasks.  
 585 *Cognitive Psychology*. 2010;60(3):127-57.
- 586 53. Jung SH. Quasi-likelihood for median regression models. *Journal of the American Statistical*  
 587 *Association*. 1996;91(433):251-7.
- 588 54. Koenker R, Bassett G. Regression Quantiles. *Econometrica*. 1978;46(1):33-50.
- 589 55. Koenker R, Hallock KF. Quantile regression. *Journal of Economic Perspectives*.  
 590 2001;15(4):143-56.
- 591 56. Wei Y, Pere A, Koenker R, He XM. Quantile regression methods for reference growth  
 592 charts. *Statistics in Medicine*. 2006;25(8):1369-82.
- 593 57. Expert Panel on Detection E, Treatment of High Blood Cholesterol in A. Executive  
 594 Summary of The Third Report of The National Cholesterol Education Program (NCEP)  
 595 Expert Panel on Detection, Evaluation, And Treatment of High Blood Cholesterol In Adults  
 596 (Adult Treatment Panel III). *JAMA*. 2001;285(19):2486-97.
- 597 58. Cox DR. Regression Models and Life-Tables. *Journal of the Royal Statistical Society Series*  
 598 *B-Statistical Methodology*. 1972;34(2):187-+.
- 599 59. Cox DR. Partial Likelihood. *Biometrika*. 1975;62(2):269-76.

- 600 60. Dempster AP, Laird NM, Rubin DB. Maximum Likelihood from Incomplete Data Via Em  
 601 Algorithm. *Journal of the Royal Statistical Society Series B-Methodological*. 1977;39(1):1-  
 602 38.
- 603 61. Laird NM, Ware JH. Random-Effects Models for Longitudinal Data. *Biometrics*.  
 604 1982;38(4):963-74.
- 605 62. Rubin DB. Inference and Missing Data. *Biometrika*. 1976;63(3):581-90.
- 606 63. Ogihara T, Miki M, Kitagawa M, Mino M. Distribution of Tocopherol among Human-  
 607 Plasma Lipoproteins. *Clinica Chimica Acta*. 1988;174(3):299-305.
- 608 64. Winbauer AN, Pingree SS, Nuttall KL. Evaluating serum alpha-tocopherol (vitamin E) in  
 609 terms of a lipid ratio. *Ann Clin Lab Sci*. 1999;29(3):185-91.
- 610 65. Vanvliet T, Schreurs WHP, Vandenberg H. Intestinal Beta-Carotene Absorption and  
 611 Cleavage in Men - Response of Beta-Carotene and Retinyl Esters in the Triglyceride-Rich  
 612 Lipoprotein Fraction after a Single Oral Dose of Beta-Carotene. *American Journal of Clinical*  
 613 *Nutrition*. 1995;62(1):110-6.
- 614 66. Redlich CA, Chung JS, Cullen MR, Blaner WS, Van Bennekum AM, Berglund L. Effect of  
 615 long-term beta-carotene and vitamin A on serum cholesterol and triglyceride levels among  
 616 participants in the Carotene and Retinol Efficacy trial (CARET) (vol 143, pg 427, 1999).  
 617 *Atherosclerosis*. 1999;145(2):423-+.
- 618 67. Johnson CL, Paulose-Ram R, Ogden CL, Carroll MD, Kruszon-Moran D, Dohrmann SM, et  
 619 al. National health and nutrition examination survey: analytic guidelines, 1999-2010. *Vital*  
 620 *Health Stat 2*. 2013(161):1-24.
- 621 68. Ioannidis JP, Loy EY, Poulton R, Chia KS. Researching genetic versus nongenetic  
 622 determinants of disease: a comparison and proposed unification. *Sci Transl Med*.  
 623 2009;1(7):7ps8.

- 624 69. Walter SD, Holford TR. Additive, Multiplicative, and Other Models for Disease Risks.  
 625 American Journal of Epidemiology. 1978;108(5):341-6.
- 626 70. Wolff MS, Engel SM, Berkowitz GS, Ye X, Silva MJ, Zhu C, et al. Prenatal phenol and  
 627 phthalate exposures and birth outcomes. Environ Health Perspect. 2008;116(8):1092-7.
- 628 71. Varshavsky JR, Zota AR, Woodruff TJ. A Novel Method for Calculating Potency-Weighted  
 629 Cumulative Phthalates Exposure with Implications for Identifying Racial/Ethnic Disparities  
 630 among U.S. Reproductive-Aged Women in NHANES 2001-2012. Environ Sci Technol.  
 631 2016;50(19):10616-24.
- 632 72. Howard GJ, Webster TF. Contrasting theories of interaction in epidemiology and toxicology.  
 633 Environ Health Perspect. 2013;121(1):1-6.
- 634 73. VanderWeele TJ. On the Distinction Between Interaction and Effect Modification.  
 635 Epidemiology. 2009;20(6):863-71.
- 636 74. Pedersen EJ, Miller DL, Simpson GL, Ross N. Hierarchical generalized additive models in  
 637 ecology: an introduction with mgcv. PeerJ. 2019;7:e6876.
- 638 75. Foster JC, Taylor JMG, Nan B. Variable selection in monotone single-index models via the  
 639 adaptive LASSO. Stat Med. 2013;32(22):3944-54.
- 640 76. Yang H, Yang J. A robust and efficient estimation and variable selection method for partially  
 641 linear single-index models. J Multivariate Anal. 2014;129:227-42.
- 642 77. Lai P, Wang QH, Lian H. Bias-corrected GEE estimation and smooth-threshold GEE  
 643 variable selection for single-index models with clustered data. J Multivariate Anal.  
 644 2012;105(1):422-32.
- 645 78. Friedman JH, Stuetzle W. Projection Pursuit Regression. Journal of the American Statistical  
 646 Association. 1981;76(376):817-23.
- 647

648 **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
<b>Identity link function</b>											
<b>Environmental factors</b>											
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)
<b>Covariates</b>											
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)		
<b>Ethnicity</b>											
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)		
<b>Quadratic link function</b>											
<b>Environmental factors</b>											
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)
<b>Covariates</b>											
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)		
<b>Ethnicity</b>											
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)		

649

650

651

652

653

654

655

656

657 **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
<b>Identity link function</b>											
<b>Environmental factors</b>											
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)
<b>Covariates</b>											
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)		
<b>Ethnicity</b>											
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)		
<b>Quadratic link function</b>											
<b>Environmental factors</b>											
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)
<b>Covariates</b>											
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)		
<b>Ethnicity</b>											
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)		

658

659

660

661

662

663

664

665 **Figure titles and legends**

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

668 **Fig. 2** Results from PLSI linear regression and multivariable linear regression in NHANES  
669 2002-2003 (d=8, N=800). Bars show the estimated relative importance (absolute value of  
670 estimated coefficient) of 8 environmental factors on continuous triglycerides. Red/green color  
671 represents positive/negative effect. Error bars indicate 95% CIs.

672 **Fig. 3** Results from PLSI quantile regression and multivariable quantile regression in NHANES  
673 2002-2003 (d=8, N=800). Bars show the estimated relative importance (absolute value of  
674 estimated coefficient) of 8 environmental factors on three quartiles (25th, 50th, and 75th  
675 percentiles) of triglycerides. Red/green color represents positive/negative effect. Error bars  
676 indicate 95% CIs.

677 **Fig. 4** Results from PLSI logistic regression and multivariable logistic regression in NHANES  
678 2002-2003 (N=800). Bars show the estimated relative importance (absolute value of estimated  
679 coefficient) of 8 environmental factors on dichotomized triglycerides. Red/green color represents  
680 positive/negative effect. Error bars indicate 95% CIs.