

Workshop on Analyzing Mixtures in Environmental Health Studies: WQS Regression

Chris Gennings

Icahn School of Medicine at Mount Sinai
Department of Environmental Medicine and Public Health

August 20, 2019



Overview of Mixtures

Concerns: high dimensionality; complex correlation patterns

- multicollinearity and
- reversal paradox
- Sensitivity and specificity identifying ‘bad actors’

Strategies:

- Reducing dimensionality: e.g., PCA
- Addressing ill-conditioning in regression with constraints
 - Shrinkage methods – e.g., LASSO
 - WQS regression
- Flexible response surface methods
 - e.g., Bayesian Kernel Machine Regression (BKMR)

Multicollinearity

- Correlation among predictor variables impact the variability of parameter estimates in regression models.
- The prediction of the model at observed data points may be adequate (i.e., “the old picket fence” analogy), but hypothesis tests of model parameters have decreased power.

Reversal paradox

Illustration: Assume $\text{Corr}(y, x_1)=0.2$ and $\text{Corr}(y, x_2)=0.1$. The beta estimates in a linear model are impacted by the $\text{Corr}(X_1, X_2)$.

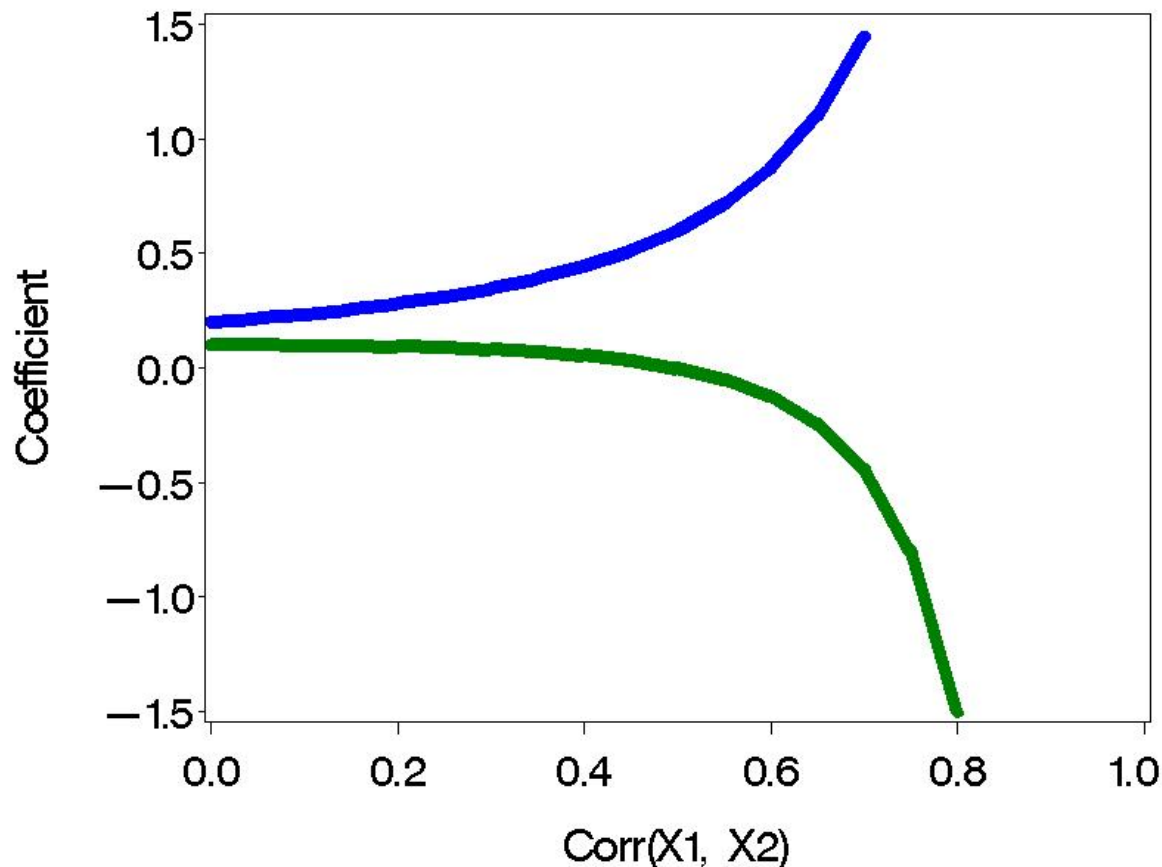


Illustration: Mitro et al, 2016 EHP

- **Background:** Exposure to persistent organic pollutants (POPs) such as **dioxins, furans, and polychlorinated biphenyls (PCBs)** may influence **leukocyte telomere length (LTL)**, a biomarker associated with chronic disease.

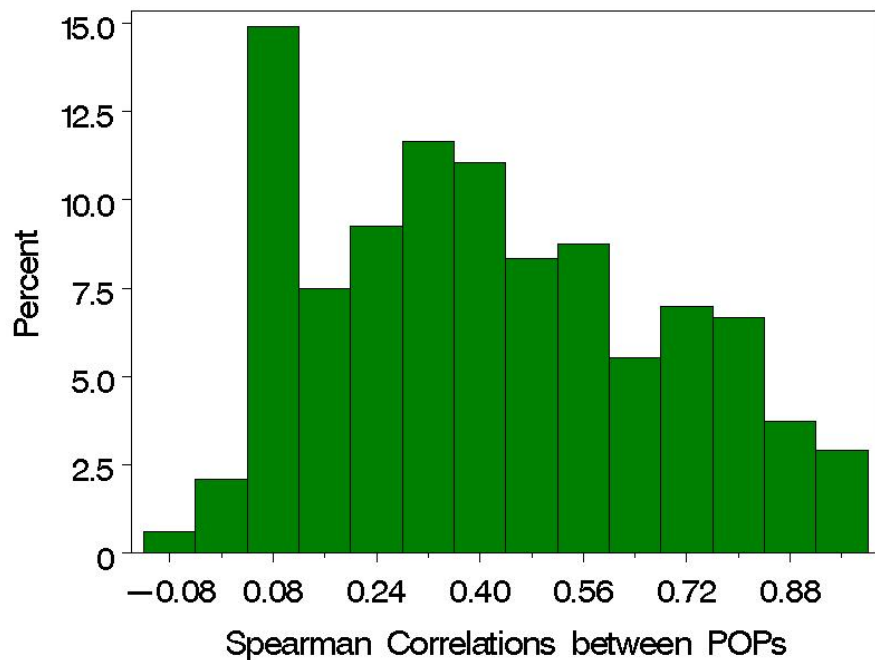
In vitro research suggests dioxins may bind to the aryl hydrocarbon receptor (AhR) and induce telomerase activity, which elongates LTL.

However, few epidemiologic studies have investigated associations between POPs and LTL.

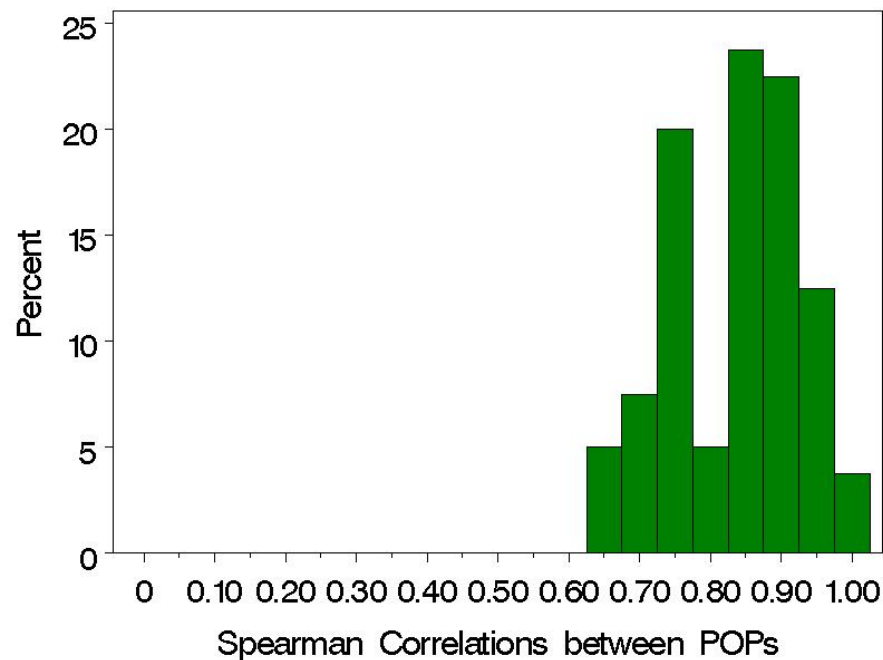
- **Covariates:**

Models were adjusted for age, age², sex, race/ethnicity, BMI, log(cotinine), white blood cell count, percent lymphocytes, percent monocytes, percent neutrophils, percent eosinophils, percent basophils

Correlation Between POPs



Full set of 18 POPs



Subset of 9 PCBs

**STABILITY OF ILL-CONDITIONING
WITH CONSTRAINTS:
VARIANCE VS BIAS**

Least Squares with Constraints

- Ridge Regression

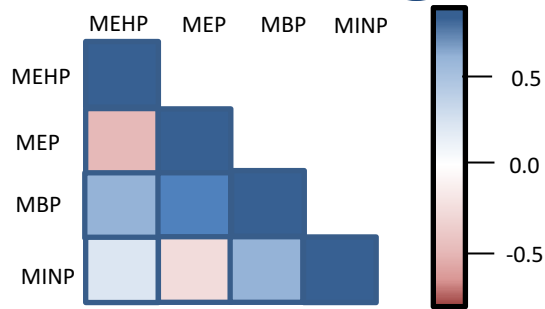
$$\hat{\beta}_{ridge} = \min_{\beta} \left[\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^{p-1} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right]$$

- LASSO

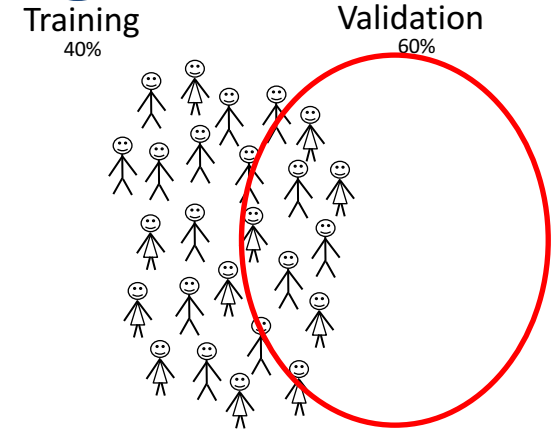
$$\hat{\beta}_{LASSO} = \min_{\beta} \left[\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^{p-1} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right]$$

- Elastic Net
- $$\hat{\beta}_{elastic\ net} = \min_{\beta} \left[\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^{p-1} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \left(\alpha |\beta_j| + (1 - \alpha) \beta_j^2 \right) \right]$$

Weighted Quantile Sum (WQS) Regression

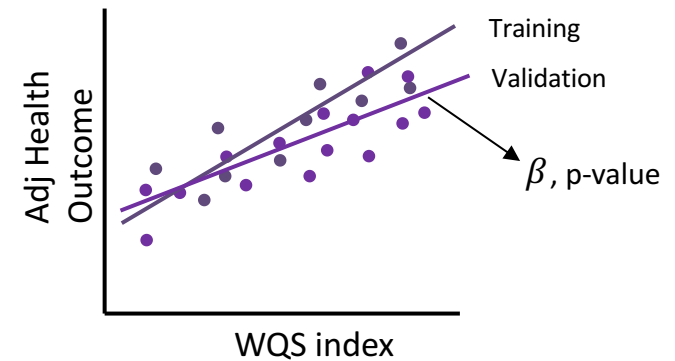
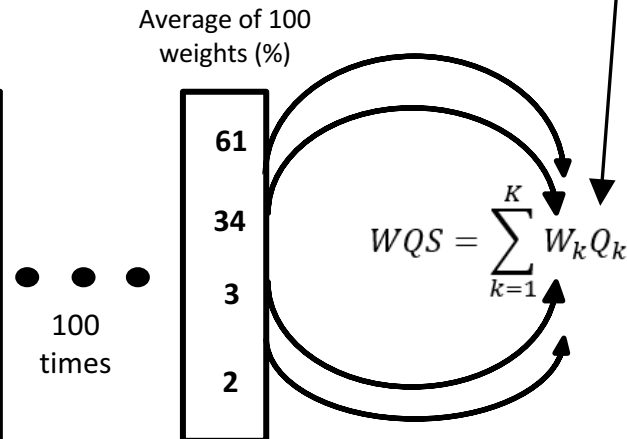


Subject ID	Concentration of MEHP (ng/ml)	Rank (Q) of MEHP
1	3.4	4
2	1.2	2
3	10.3	9



Weights (W, %)

MEHP	62	58	67	59
MEP	34	36	31	35
MBP	3	3	1	2
MINP	1	3	1	4



Weighted Quantile Sum (WQS)

Regression (Carrico et al, 2014)

Nonlinear regression with weight parameters:

$$\theta = [\beta_0, \beta_1, w_1, \dots, w_c, \gamma']$$

$$g(\mu) = \beta_0 + \beta_1 \sum_{j=1}^c w_j q_j + \sum_{k=1} \gamma_k z_{ik}$$

Final WQS index is a weighted average across the bootstrap samples using a 'signal function'

$$WQS = \sum_{j=1}^c \bar{w}_j q_j$$

$$\bar{w}_j = \frac{1}{B} \sum_{b=1}^B w_{j(b)} f(\hat{\beta}_{1(b)})$$

Final model:

$$g(\mu) = \beta_0 + \beta_1 WQS + \sum_{k=1} \gamma_k z_{ik}$$

Weighted Quantile Sum (WQS)

Regression (Carrico et al, 2014)

Nonlinear regression with weight parameters:

$$\theta = [\beta_0, \beta_1, w_1, \dots, w_c, \gamma']$$

$$g(\mu) = \beta_0 + \beta_1 \sum_{j=1}^c w_j q_j + \sum_{k=1} \gamma_k z_{ik}$$

Why quantiles?

Final WQS index is a weighted average across the bootstrap samples using a 'signal function'

$$WQS = \sum_{j=1}^c \bar{w}_j q_j$$

$$\bar{w}_j = \frac{1}{B} \sum_{b=1}^B w_{j(b)} f(\hat{\beta}_{1(b)})$$

Final model:


$$g(\mu) = \beta_0 + \beta_1 WQS + \sum_{k=1} \gamma_k z_{ik}$$

Nonlinear Least Squares with Constraints

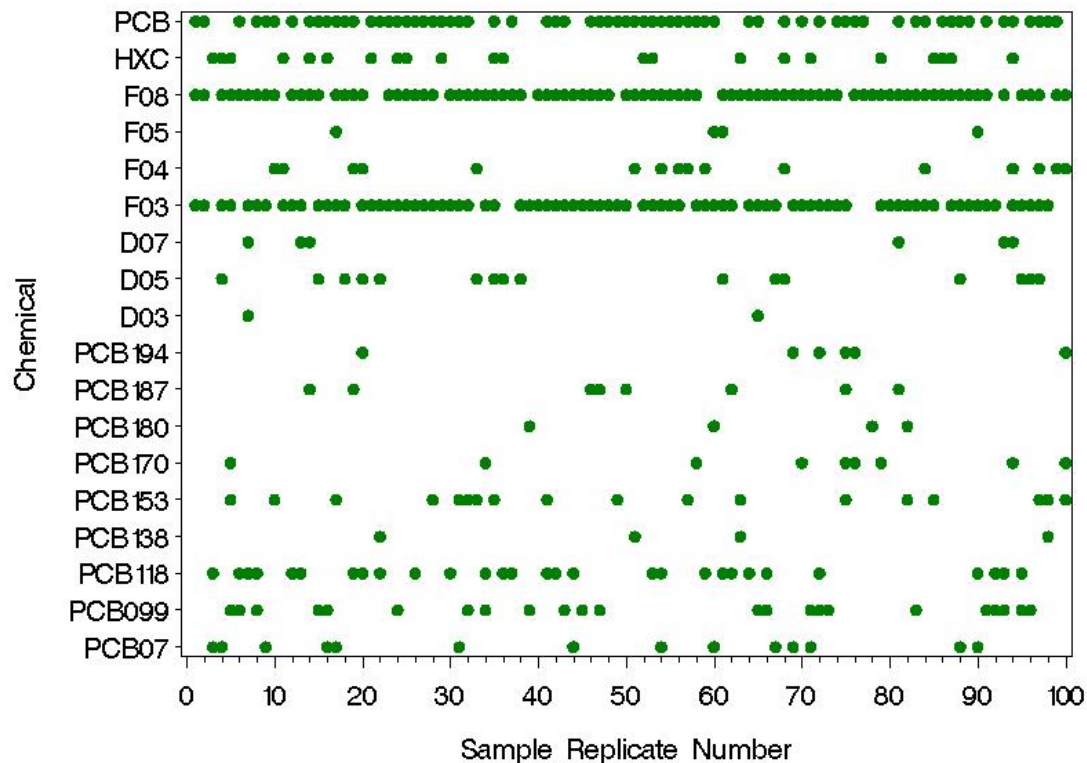
- WQS Regression with a Lagrange multiplier and with implicit directionality constraint

$$\hat{\theta}_{WQS} = \min_{\beta} \left[\sum_{i=1}^n \left(y_i - \left(\beta_0 + \beta_1 \sum_{j=1}^c w_j q_j + \sum_{k=1} \gamma_k z_{ik} \right) \right)^2 + \lambda \left(\sum_{j=1}^c w_j - 1 \right) \right]$$

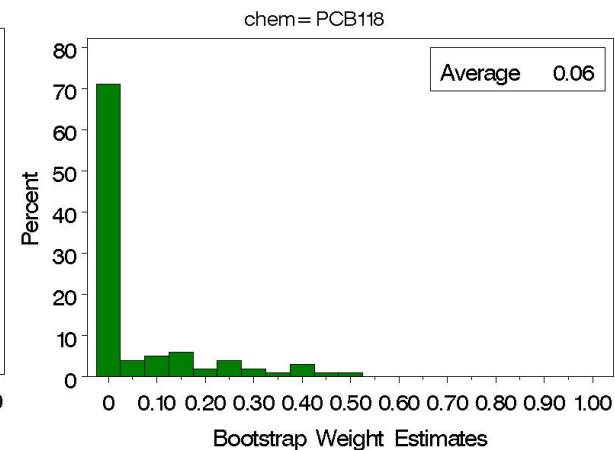
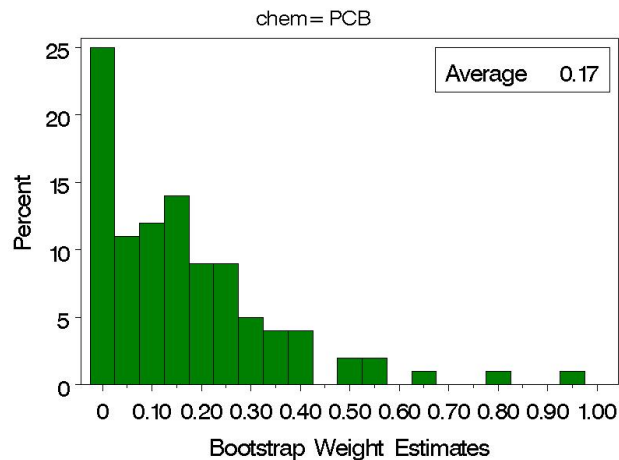
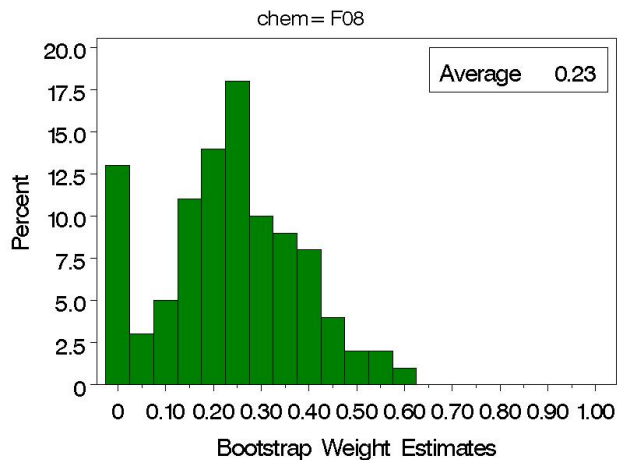
WQS regression: Ensemble step

- Bootstrap samples of *observations*
 - Why?
 - How many samples?
 - Distribution of weights
 - Random subset of *components* (i.e., c variables)
 - Subsets of size, say, \sqrt{c}
 - 1000 random subsets
 - Average across full set
- 
- Two Strategies

Distribution of Weights across Bootstrap Samples



wtGT05 0 1



Splitting data for Training & Testing

- Generally, we use 40% of the sample for estimating weights and 60% for testing significance of the index
- Need more power for testing for significance of β_1
- Impact of random split:
 - Eva Tanner is working on using repeated random holdouts (Tanner et al, submitted)



EXAMPLE: 9 PCBs and LTL

Preliminary adjusted analyses

Single chemical

Parameter	Estimate	StdErr	ProbChiSq
log_LBX074LA	0.128	0.022	13E-9
log_LBX099LA	0.107	0.022	62E-8
log_LBX118LA	0.112	0.019	8E-9
log_LBX138LA	0.097	0.02	16E-7
log_LBX153LA	0.104	0.021	12E-7
log_LBX170LA	0.094	0.026	33E-5
log_LBX180LA	0.073	0.023	0.001
log_LBX187LA	0.085	0.024	46E-5
log_LBX194LA	0.061	0.028	0.032

Joint model

Parameter	Estimate	Standard Error	Pr > ChiSq
logLBX074LA	0.0339	0.0197	0.0849
logLBX099LA	0.0037	0.0221	0.8674
logLBX118LA	0.0087	0.0193	0.6543
logLBX138LA	-0.0360	0.0354	0.3095
logLBX153LA	0.0904	0.0421	0.0315
logLBX170LA	-0.0015	0.0368	0.9664
logLBX180LA	-0.0348	0.0283	0.2181
logLBX187LA	-0.0077	0.0253	0.7603
logLBX194LA	-0.0019	0.0264	0.9423

EXAMPLE: WQS regression

Split: 40% for estimating weights; 60% for testing significance of WQS index

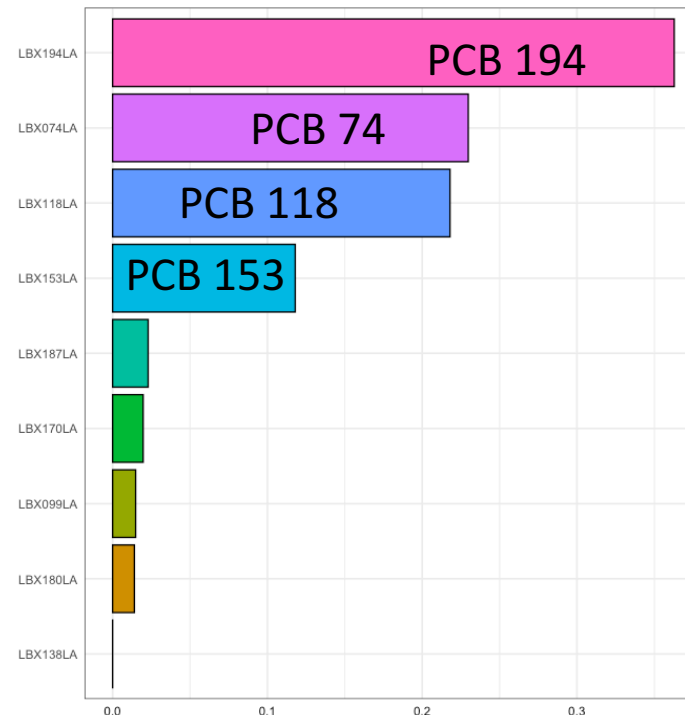
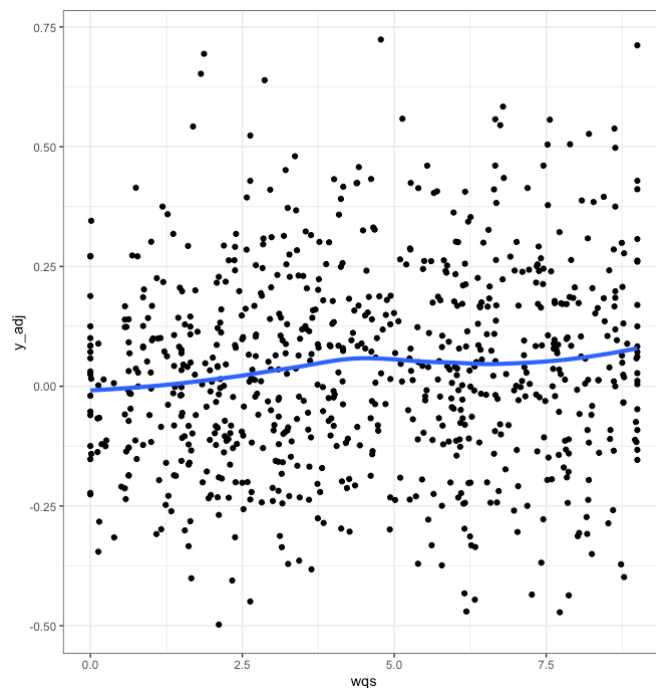
Quantiles: **deciles**

100 bootstrap samples

Analysis adjusted by covariates

Beta1 unconstrained

Cut-point for identifying a “bad actor”: $1/9 = 0.11$



Beta1 = 0.021
SE= 0.005
p <0.001

EXAMPLE: WQS Regression

Split: 40% for estimating weights; 60% for testing significance of WQS index

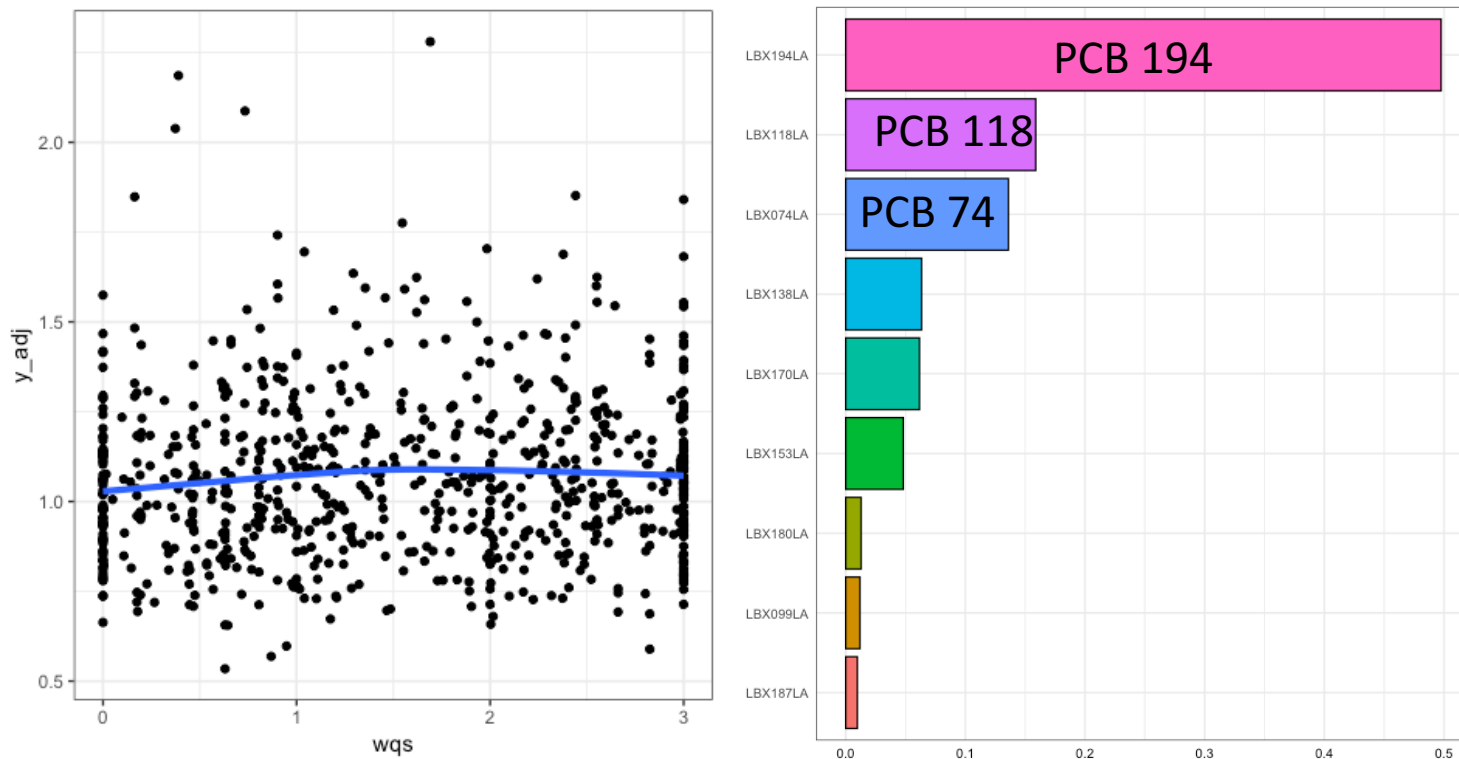
Quantiles: **quartiles**

100 bootstrap samples

Analysis adjusted by covariates

Beta1 unconstrained

Cut-point for identifying “bad actor”: $1/9 = 0.11$



Beta1 = 0.040
SE= 0.012
P = 0.001

Stratified WQS regression

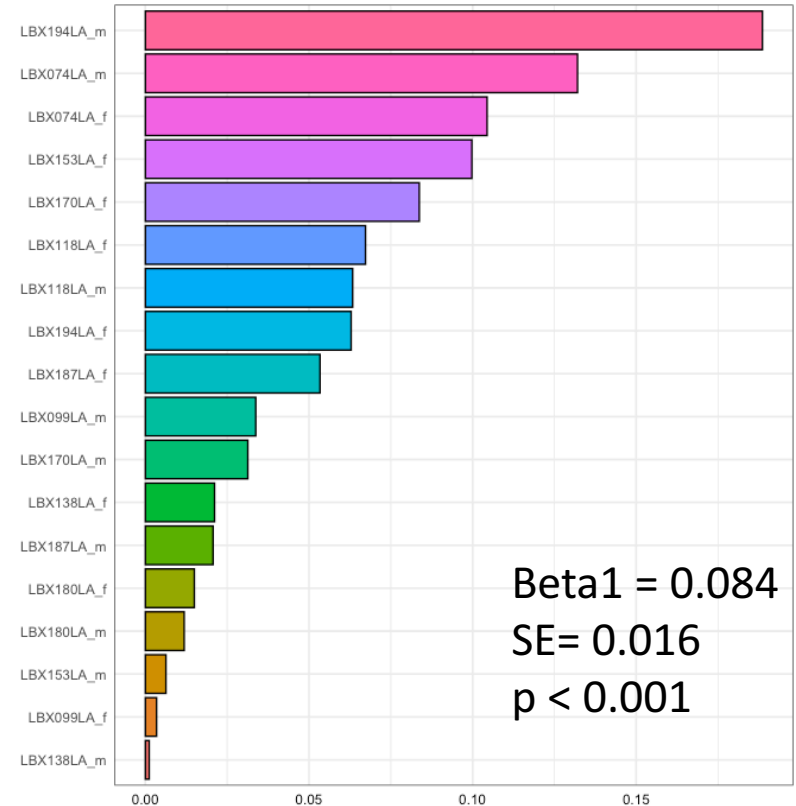
(Brunst et al, 2017, AJE)

- Similar to interaction between a categorical variable and the weights
- Weights are estimated per each category in a single index where weights sum to 1
- STEPS:
 - Determine overall quantiles per component
 - Use interaction quantile scoring; e.g.,

$$qchem_{males} = \begin{cases} q, & \text{if male} \\ 0, & \text{otherwise} \end{cases} \cdot qchem_{females} = \begin{cases} q, & \text{if female} \\ 0, & \text{otherwise} \end{cases}$$

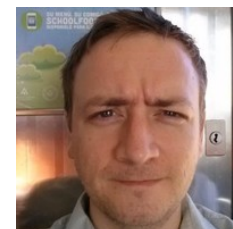
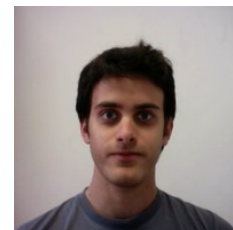
Stratified WQS regression

PCB	Male	Rel Wt(%)	Female	Rel Wt(%)
74	0.13	27	0.10	20
99	0.03	6	<0.01	1
118	0.06	12	0.07	14
138	<0.01	0	0.02	4
153	0.01	2	0.10	20
170	0.03	6	0.08	16
180	0.01	2	0.01	2
187	0.02	4	0.05	10
194	0.19	40	0.06	12
SUM	0.49		0.51	



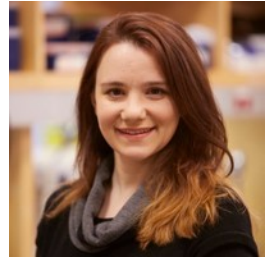
Wrap-up

- Ill-conditioning due to multicollinearity in environmental health data is improved by constraints in the optimization for parameter estimation.
- Choice of strategy depends on the research question:
 - Biomarker identification (e.g., shrinkage methods)
 - Mixture effect (e.g., PCA, WQSR, BKMR)
 - Interaction among components (e.g., BKMR)
- **WQS regression** is based on quantile scores and is improved with the addition of the ensemble step
 - It addresses questions of a mixture effect with an empirically weighted index;
 - Stratified WQSR has the advantage that the sample size is not reduced to each strata
 - R packages: gWQS and WQSrs are being developed by **Stefano Renzetti** and **Paul Curtin**
 - Extensions are forthcoming...
 - Repeated hold-outs is being studied by **Eva Tanner**
 - A Bayesian version is being developed by **Elena Colicino**



Acknowledgments

- **A team is what we have....**
 - Elena Colicino
 - Paul Curtin
 - Stefano Renzetti
 - Eva Tanner
- **Funding sources:**
 - NIH (T32 ES007334; U2CES026555 ; R01ES028811)

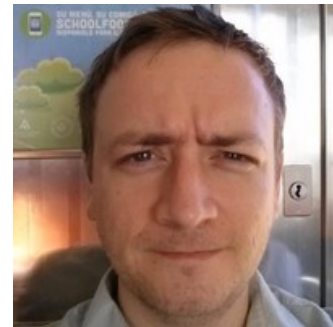


THANK YOU!

EXTENSIONS AND RECENT WORK

Extensions to R package

Stefano Renzetti is the developer of the gWQS and WQSrs R packages with assistance by Paul Curtin



Data type extensions

- New gWQS package includes capability for evaluating new data types/distributions:
 - Multinomial (generalized logits)
 - Poisson (count data)
 - Negative binomial (over-dispersed count data)
 - Stratification implementation for categorical variables
- Later extensions will include
 - Time-to-response data with censoring (e.g., Weibull distribution)
 - Allowance for interaction of WQS with continuous variable during estimation

Random subset WQS regression

- Two types of ensemble steps
 - *Bootstrap sampling* of observations with replacement
 - *Random subset* selection of variables (e.g., random set of 20 repeated 1000 times)

Allows for WQS regression to be extended to large number of variables – e.g., metabolomics

Extensions: Metabolomics

The methods of metabolomics are not only to understand traditional measures of **biological response** but also to analyze the **exposures** associated with those responses.

May be useful for

- Biomarker discovery
- Measuring a “mixture effect”

RS WQS regression seems to work well in high dimensions (Curtin et al 2019, *Comm in Stats*)



Repeated Hold-outs WQS regression

- **Eva Tanner** has extended WQS regression to include repeated hold-outs to accommodate the variability due to the random seed in splitting the data (Tanner et al, under revision review)

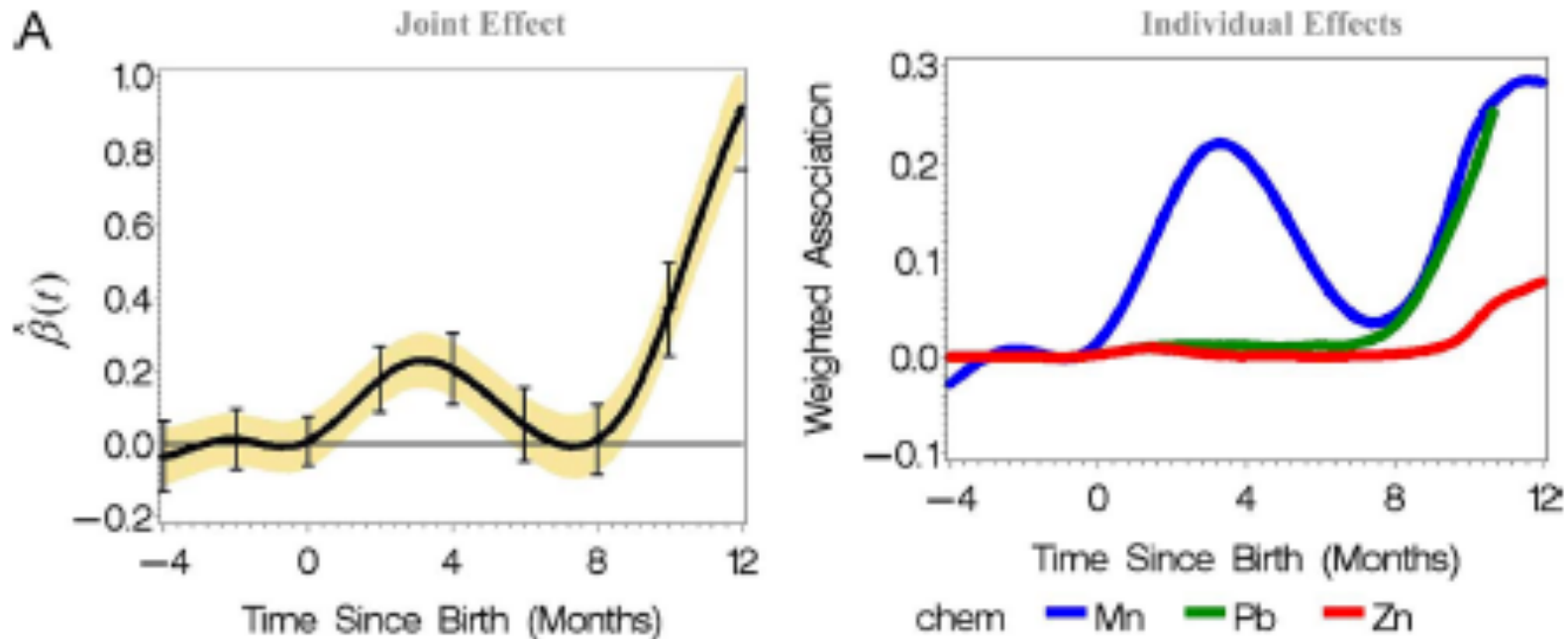


Lagged WQS Regression

(Bello et al, Env Res, 2017)

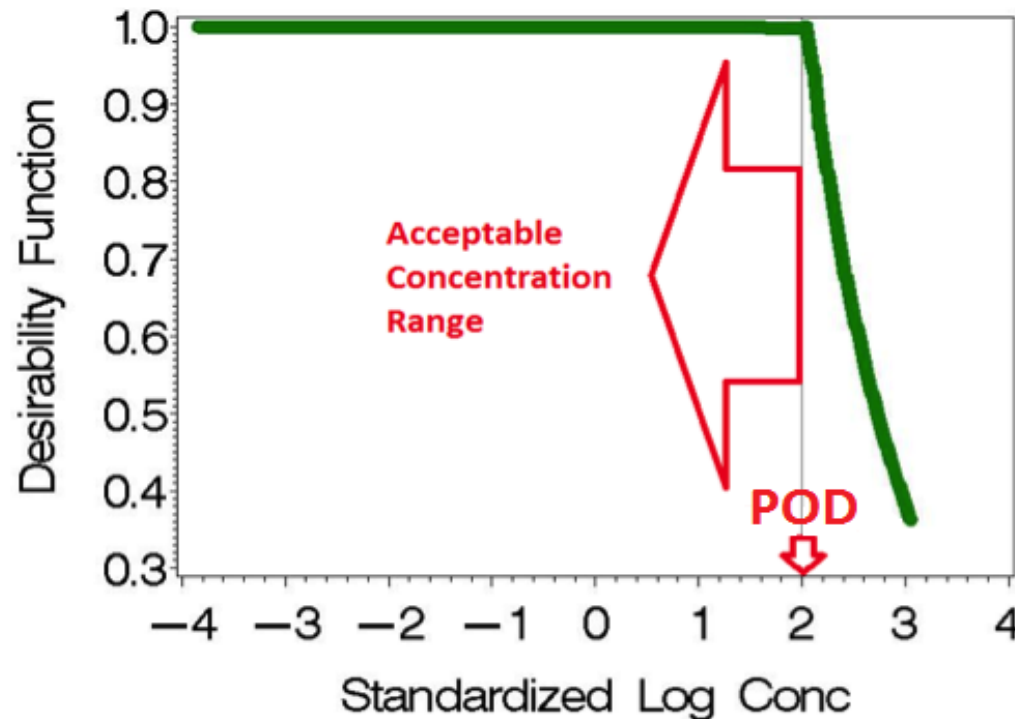
revised algorithm – in the works...

Lagged WQS regression is a reverse DLM on an iteratively weighted WQS index.



Acceptable Concentration Region (ACR) models (Gennings et al 2018 ENV INT)

Incorporates the concept of **regulatory guideline values** into a nonlinear regression model



ACR model example

For single chemicals

$$d_m^{low} = \begin{cases} 1, & X_m < \delta_m^{low} \\ \exp[-\gamma_m^{low} (X_m - \delta_m^{low})], & X_m \geq \delta_m^{low} \end{cases}$$

$$g(\mu_i) = \begin{cases} \beta_0 + \beta_1(1) + Z_i^T \theta, & X_i < \delta^{low} \\ \beta_0 + \beta_1(\exp[-\gamma^{low} (X_i - \delta^{low})]) + Z_i^T \theta, & X_i \geq \delta^{low} \end{cases}$$

For mixtures $g(\mu_i) = \beta_0 + \beta_1(d_1 \times d_2 \times \dots \times d_M)^{\frac{1}{M}} + Z_i^T \gamma$

$$= \beta_0 + \beta_1 MDF + Z_i^T \gamma$$

Extensions are underway with Eva Turner

THANK YOU!