

Repeated Holdout Validation for WQS

Eva M Tanner, PhD, MPH

Postdoctoral Researcher

Icahn School of Medicine at Mount Sinai

etanner@mssm.edu

<https://github.com/evamtanner>

Mentor: Chris Gennings

ISEE 2019 Workshop

Mixtures Analysis with Weighted Quantile Sum (WQS) Regression and its Extensions

8/25/19



**Mount
Sinai**

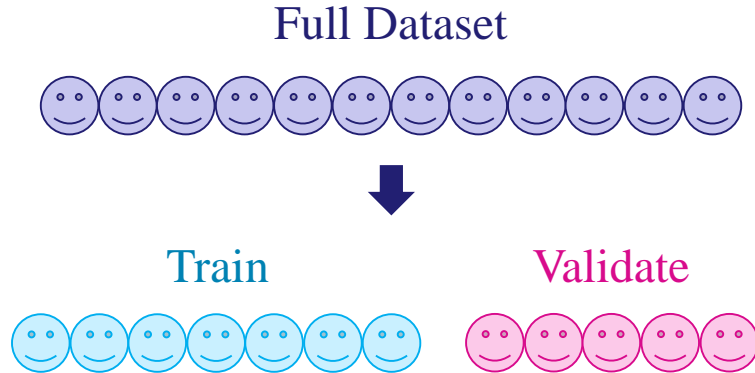
Outline

- ▶ **Cross-validation Techniques**
- ▶ **Repeated Holdout Validation**
- ▶ **Application to SELMA Study**
- ▶ **R Tutorial**

Cross-validation Techniques

What is Cross-Validation?

- ▶ Used in **predictive modeling** & machine learning for variable/model selection & to **evaluate model performance** (replicability of results)



- ▶ Can also be used in explanatory (hypothesis driven) modeling to avoid fitting to noise & to assess generalizability

Some Types of Validation

► Single Split

Partition

1



Train

Validate

Final Estimate

β & p-value

► K-fold

1



2



3



► Leave-One-Out (LOO)

1



2



3



4



⋮

12



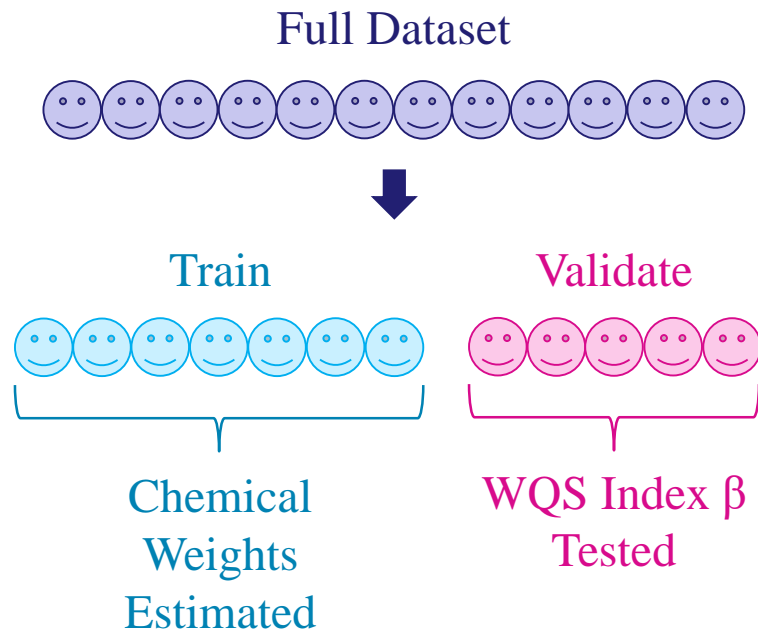
typically
k = 5-10

$$\bar{\beta} = \frac{1}{k} \sum_{i=1}^k \beta_i$$

n = sample size

$$\bar{\beta} = \frac{1}{n} \sum_{i=1}^n \beta_i$$

Prior WQS Applications used a Single Split



$$Y = \beta_0 + \beta_1 \left(\sum_{i=1}^c w_i q_i \right) + \text{covariates}$$

Outcome Y

Intercept β_0

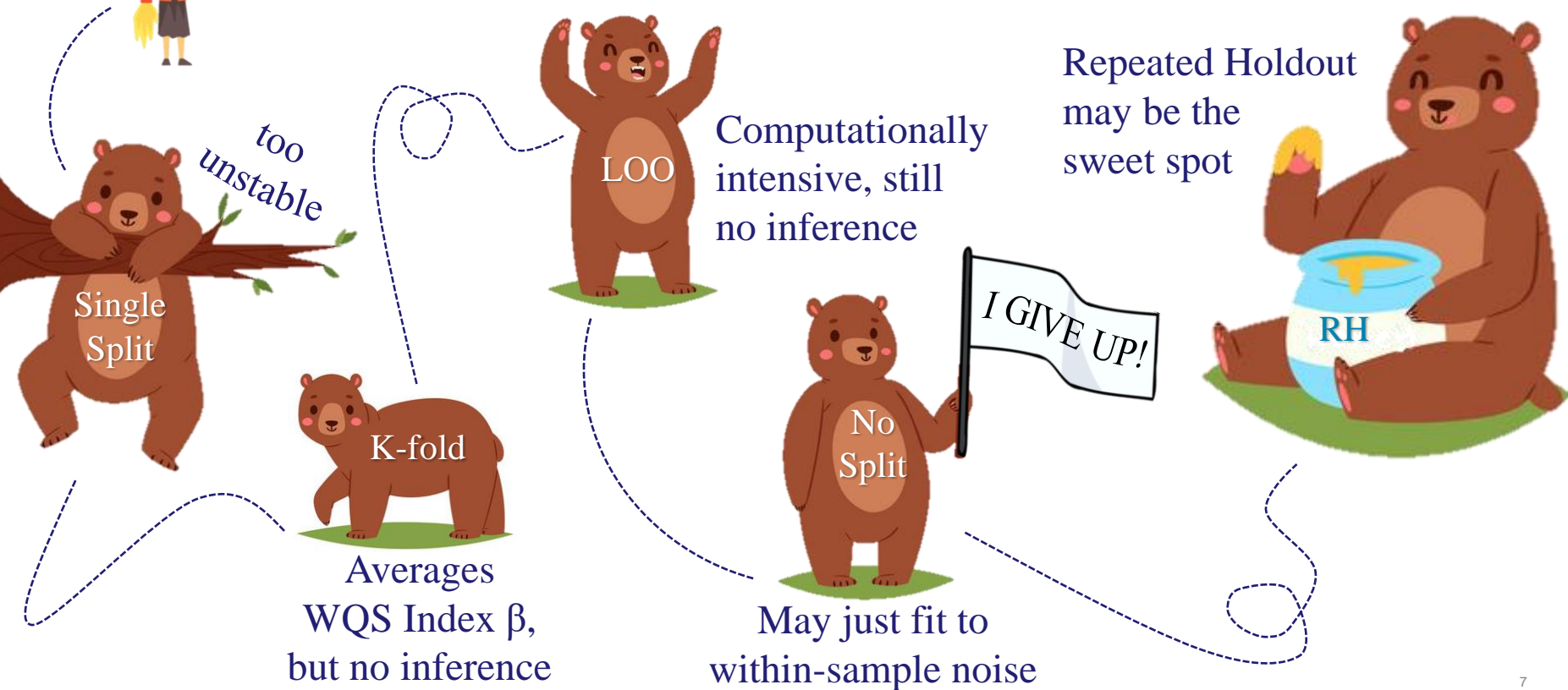
WQS Index Estimate β_1

c chemicals binned into quantile q_i with weight $0 \leq w_i \leq 1$

In smaller sample sizes
a single split can lead to
unrepresentative partitions
and **unstable estimates**



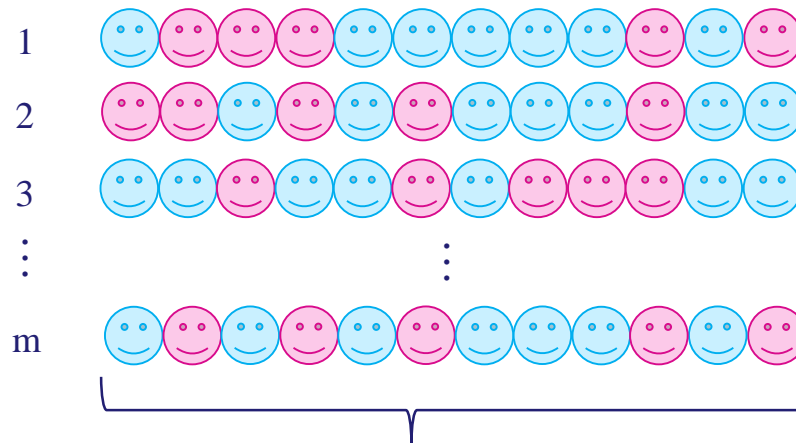
Which validation technique is just right for WQS applied to smaller (epi-relevant) sample sizes?



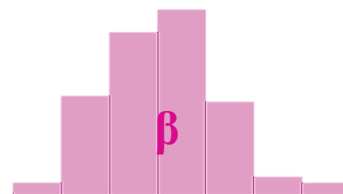
Repeated Holdout Validation

- ▶ Single Split + m Bootstraps
- ▶ Enables nonparametric inference of WQS Index estimate
 - ▶ β_{50} ($\beta_{2.5}$, $\beta_{97.5}$)
- ▶ Characterize uncertainty in selecting chemicals of concern

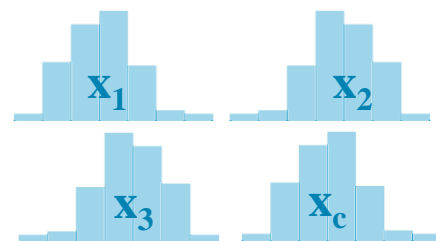
Partition



Distribution of validated
WQS Index estimates



Distribution of weights
for c chemicals



Application to SELMA Study

Endocrine Disrupting Chemicals (EDCs) = Xenobiotics able to Interfere with Hormone Action



PHENOLS



**PERSISTENT
ORGANIC
POLLUTANTS
(POPs)**



PLASTICIZERS



PESTICIDES

Prenatal EDC Exposure Impacts Child Neurodevelopment

- ▶ POPs, organophosphate & pyrethroid pesticides, phthalates, & BPA associated with
 - ▶ Altered infant brain development
 - ▶ Lower cognitive functioning
 - ▶ Neurobehavioral changes
- ▶ LIMITATION: Single chemicals evaluated in isolation
- ▶ GOAL: Evaluate impact of EDCs mixture on child IQ

Swedish Environmental Longitudinal Mother and Child, Asthma and Allergy (SELMA)

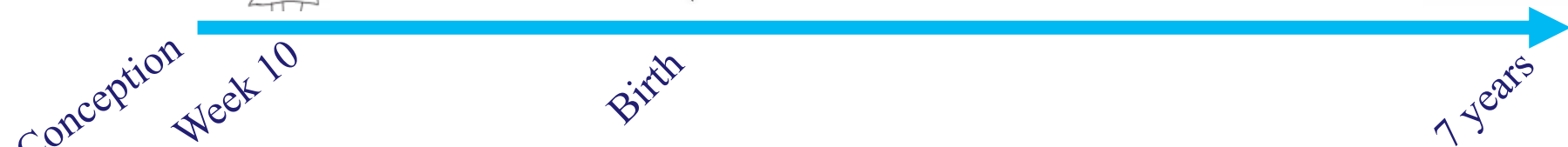
EDCs measured in
blood & urine
(N=2325)



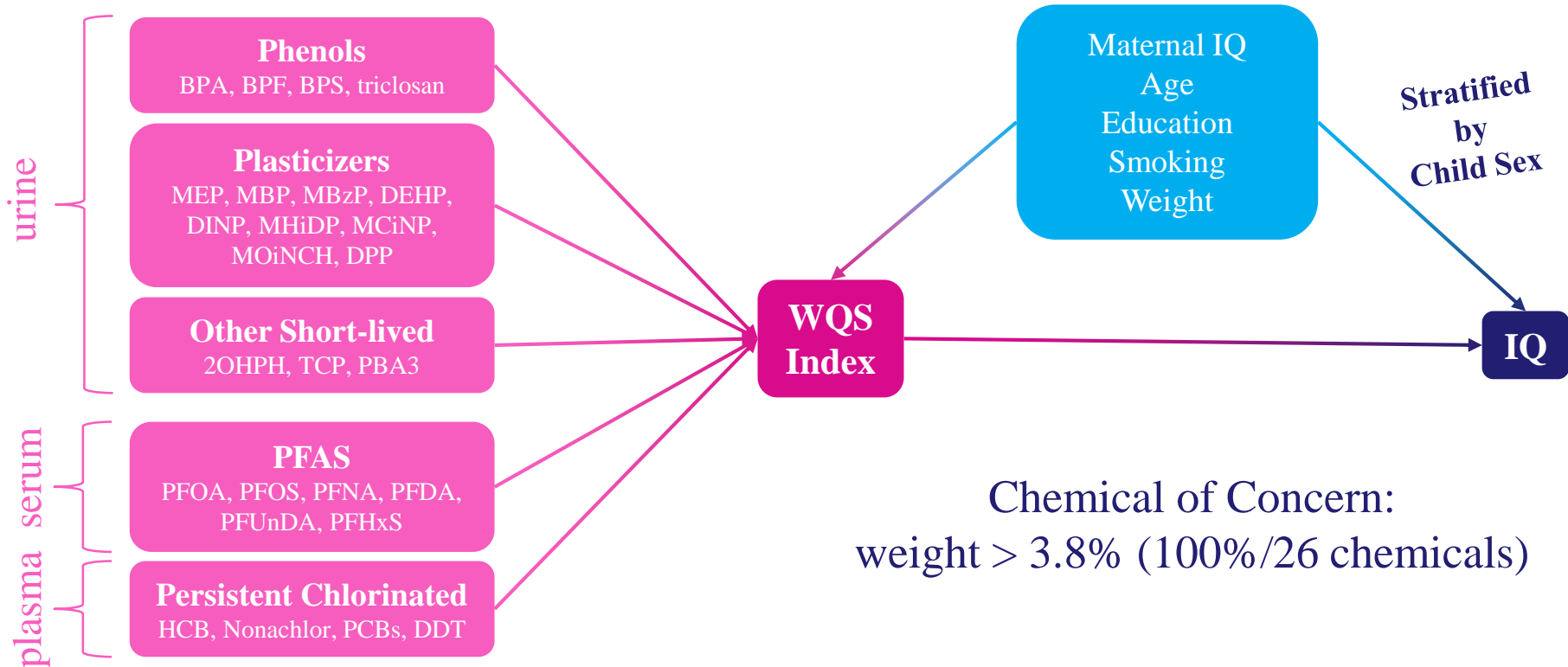
Phthalates associated
with Language Delay
(N=963)



Cognitive assessment
using WISC-IV
(N=718)



Evaluate EDCs in Relation to IQ using WQS Regression

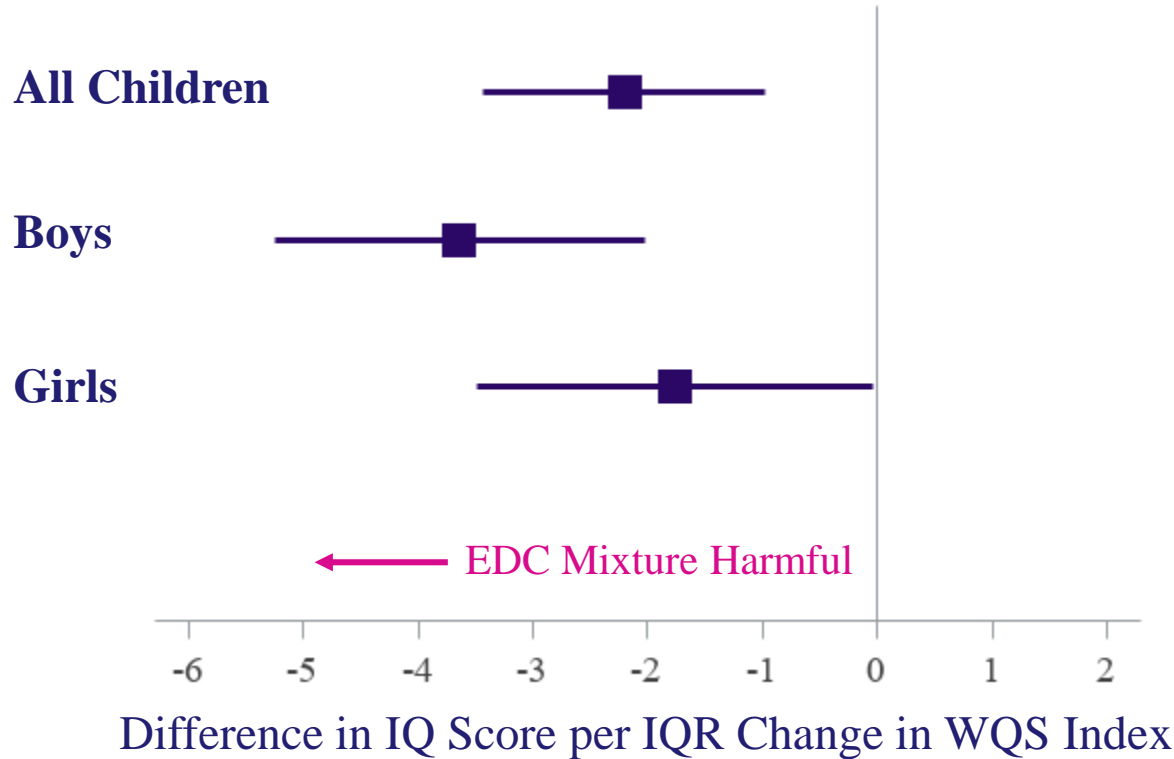


Analysis Challenges in SELMA

- ▶ Single Split
 - ▶ WQS Index β and selected chemicals changed depending on random seed
- ▶ No Split (training/testing on same data)
 - ▶ Stable estimates, but lacked rigor of validation step
- ▶ Repeated holdout a viable solution

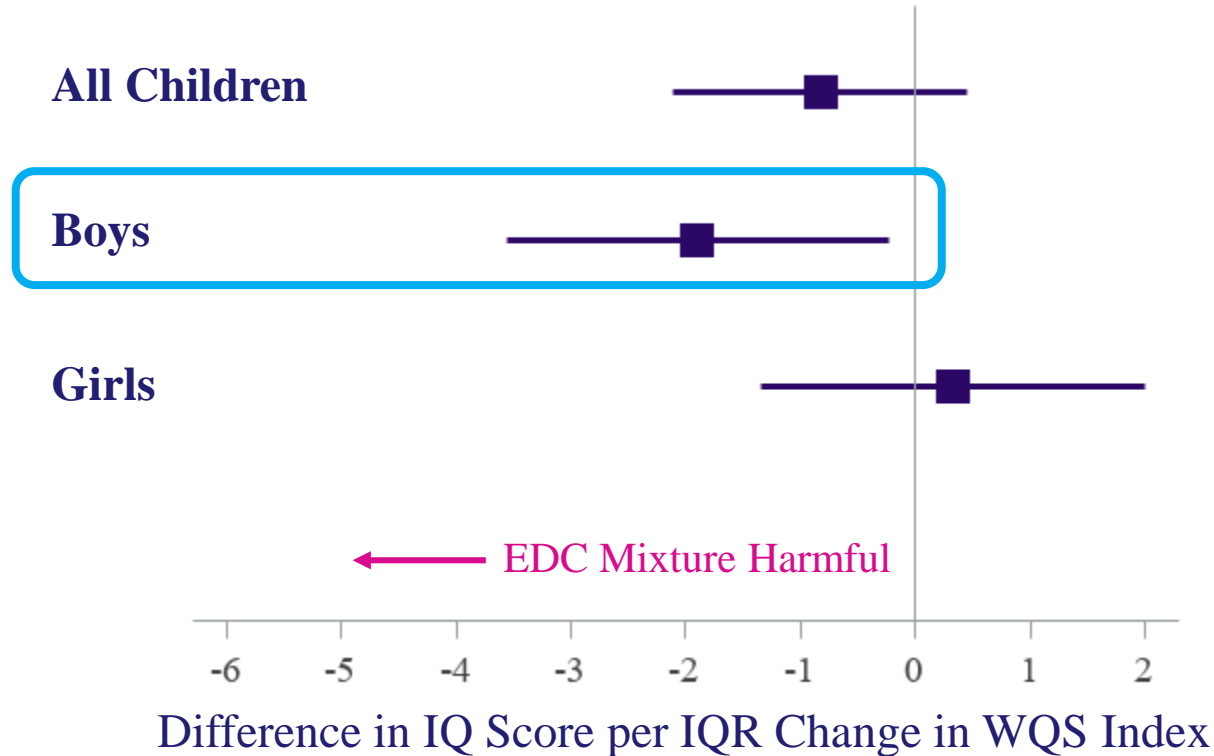
Impact of Prenatal EDC Mixture on IQ at Age 7

No Split

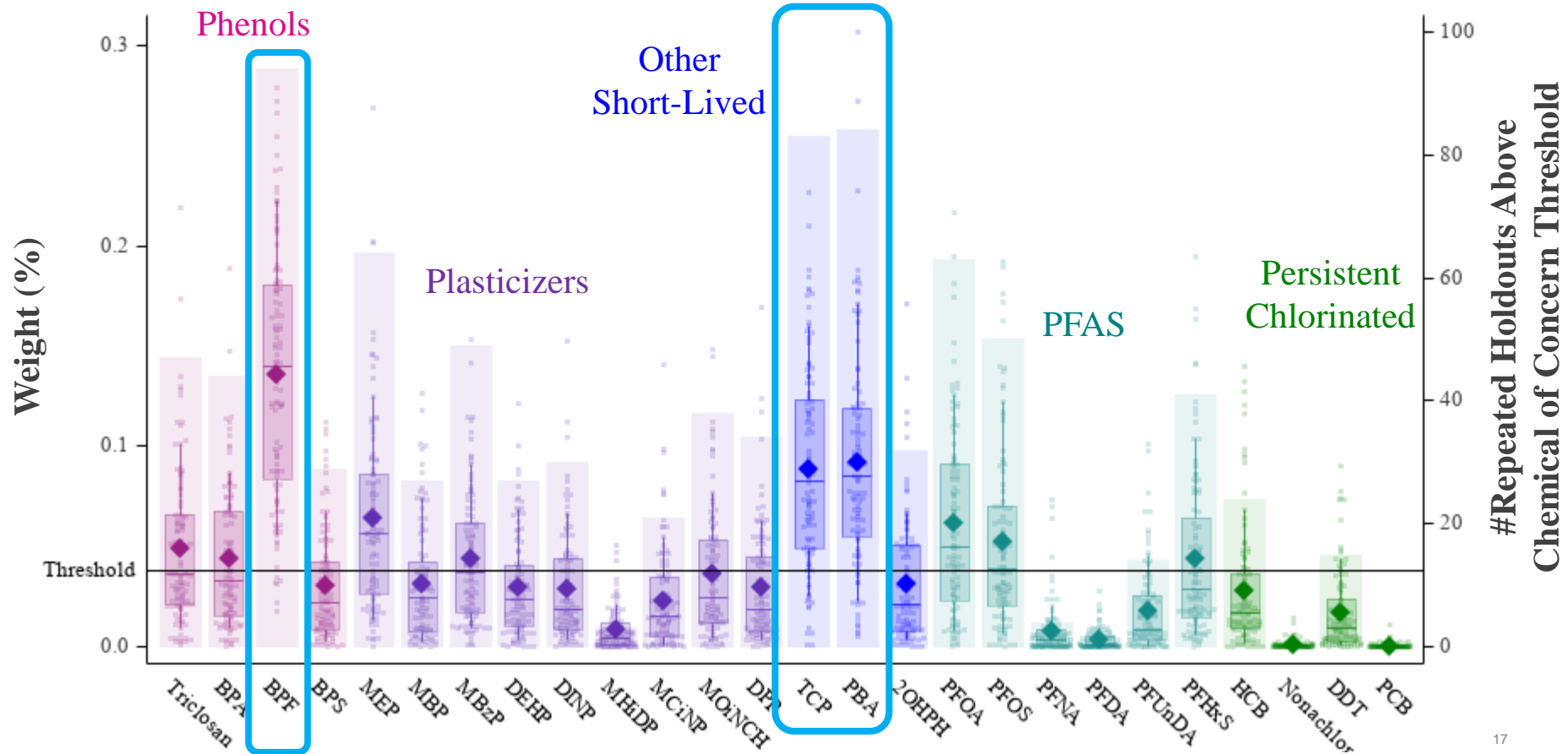


Impact of Prenatal EDC Mixture on IQ at Age 7

Repeated Holdout Validation



Chemicals of Concern Identification & Uncertainty: Boys



Summary

- ▶ WQS with single training-validation splits may lead to unrepresentative partitions & unstable results in finite samples
 - ▶ Test this by rerunning WQS with different random seed
- ▶ Training/testing on same data not necessarily wrong, but may reflect within-sample noise & results may not generalize
- ▶ Repeated Holdout Validation applied to WQS allows
 - ▶ Inference of WQS Index β
 - ▶ Characterizes weight uncertainty
- ▶ Number holdouts required depends on number needed to approximate a $\sim N$ sampling distribution in that sample

Questions?

Acknowledgements



- ▶ **Chris Gennings**
- ▶ **Carl-Gustaf Bornehag**
- ▶ **Co-authors:**
 - ▶ Maria Unenge Hallerbäck
 - ▶ Sverre Wikström
 - ▶ Christian Lindh
 - ▶ Hannu Kiviranta
- ▶ **Colleagues:**
 - ▶ Jonathan Heiss
 - ▶ Anu Joshi
- ▶ **SELMA participants**



- ▶ **Funding:**
 - ▶ EDC-MixRisk (#634880) European Union's Horizon 2020 Research and Innovation Programme
 - ▶ NIEHS PRIME: #R01ES028811-01



EDC-MixRisk
safe chemicals for future generations

**Repeated Holdout for WQS
R Tutorial
Available HERE**