# Supplemental Information for *Effect and Pathways Modifiers in a Bayesian Pathways Analysis of the National Human Exposure Assessment Survey for Arsenic in EPA Region 5*

## Definitions of Input Variables and Fitted Bayesian Pathways Models

This document specifies the variables and results from all Bayesian Pathways Models fitted to the NHEXAS data in the paper *Effect and Pathways Modifiers in a Bayesian Pathways Analysis of the National Human Exposure Assessment Survey for Arsenic in EPA Region 5*. It also lists the full-conditional distributions of the model parameters used by our MCMC algorithm.

## 1 Media and Pathway Modifier Variables

See Tables 1 and 2.

## 2 Full Conditional Distributions of the Hierarchical Bayesian Subpopulation Model

### 2.1 Posterior distribution

Letting $\boldsymbol{\theta} = (\boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\omega})$, the posterior distribution is

$$
\begin{aligned}
[\boldsymbol{\theta}, \boldsymbol{X} | \boldsymbol{Y}] \quad &\propto \quad [\boldsymbol{Y} | \boldsymbol{X}, \boldsymbol{\theta}] \, [\boldsymbol{X} | \boldsymbol{\theta}] \, [\boldsymbol{\theta}] \\
&= \quad [\boldsymbol{Y} | \boldsymbol{X}, \boldsymbol{\omega}] \, [\boldsymbol{X} | \boldsymbol{\beta}, \boldsymbol{\tau}] \, [\boldsymbol{\beta}] \, [\boldsymbol{\omega}] \, [\boldsymbol{\tau}],
\end{aligned}
$$

assuming independence of the three components within $\boldsymbol{\theta}$. The distribution is not available in closed form and so we use a Gibbs sampler Markov chain Monte Carlo (MCMC), based on an augmented (including censored and missing data) joint distribution to sample from this posterior distribution.

| Variable | Type | Definition and Units |
|---|---|---|
| UR | M | natural log of arsenic concentration in Urine ($\mu g/\ell$) |
| FD | M | natural log of arsenic intake in Food ($\mu g/day$) |
| BV | M | natural log of arsenic intake in Beverage ($\mu g/day$) |
| IA | M | natural log of arsenic concentration in Indoor Air ($ng/m^3$) |
| OA | M | natural log of arsenic concentration in Outdoor Air ($ng/m^3$) |
| PA | M | natural log of arsenic concentration in Personal Air ($ng/m^3$) |
| SL | M | natural log of arsenic concentration in Soil ($ng/g$) |
| SD | M | natural log of arsenic loading in Surface (Sill) Dust ($\mu g/cm^2$) |
| WF | m | natural log of arsenic concentration in Tap Water, Flushed ($\mu g/\ell$) |
| logcrtst | C | natural log of creatinine concentration ($\mu g/\ell$) |
| *mlogcrtn* | C | negative of the natural log of creatinine concentration ($\mu g/\ell$) |
| *logcreatinine* | C | the natural log of creatinine concentration ($\mu g/\ell$) |
| *male* | C, E | binary (0/1) variable indicating the participant is male |
| *female* | E | binary (0/1) variable indicating the participant is female |
| *child4to11/Age1* | C, E | binary (0/1) variable indicating the participant is between 4 and 11 years old |
| *age* | C, E | participant's age (years) |
| *AgeAtoB* | C, E | binary (0/1) variable indicating the participant's age is between $A$ and $B$ years old |
| *Age2* | C,E | binary (0/1) variable indicating the participant's age is between 12 and 19 years |
| *hhold* | E, P | household size |
| *log hhold* | E, P | natural log of household size |
| *minushhold* | E, P | negative of household size |
| *hrs.home* | E | average number of hours participant spends inside home/day |
| *glass* | E | average number of glasses of water participant drinks/day |
| *tap.c* | P | binary (0/1) variable indicating participant cooks with tapwater |
| *fruit* | P | binary (0/1) variable indicating participant eats fruit at least 3 days/three-months |
| *tap.d/tapwater* | E | binary (0/1) variable indicating participant drinks tapwater |
| *c.air* | E | binary (0/1) variable indicating participant's home has central air |
| *no.c.air* | E | binary (0/1) variable indicating participant's home does *not* have central air |
| *gasequip* | P | binary (0/1) variable indicating participant used gas equipment in the past week |

Table 1: For each variable used in the Bayesian hierarchical analysis, the variable symbol, its use- depending on the fitted model (M = media, C = component of the creatinine adjustment, E = effect modifier, and P = pathway modifier), and its definition, including units.

| Variable | Type | Definition and Units |
|----------|------|----------------------|
| *tobacco/tobacco1* | P | binary (0/1) variable indicating the participant is a current smoker |
| *tobacco2* | P | binary (0/1) variable indicating the participant is a former smoker |
| *workshop* | P | average number of minutes/day in an enclosed workshop |
| *work20* | P | binary (0/1) variable indicating participant works $\geq$ 20 min./day in a workshop |
| *fish* | P | binary (0/1) variable indicating participant eats fish |
| *new.fish* | P | binary (0/1) variable indicating participant eats fish more than once/month |
| *no.fish* | P | binary (0/1) variable indicating participant eats fish less than once/month |
| *logBMI* | C | Natural log of Body Mass Index $(kg/m^2)$ |
| *race3_2* | C, E | a binary(0/1) variable indicating whether the participant is African-American |
| *race3_3* | C, E | binary(0/1) variable indicating the participant's race is *not* Caucasian/African-American |
| *no.wellwater* | P | binary(0/1) variable indicating the source of running water is well water |
| *wellwater* | P | binary(0/1) variable indicating the source of running water is *not* well water |

Table 2: For each variable used in the Bayesian hierarchical analysis, the variable symbol, its use-depending on the fitted model (M = media, C = component of the creatinine adjustment, E = effect modifier, and P = pathway modifier), and its definition, including units.

## 2.2 Full conditional distributions

1. For each $i$ and $j$ we set/draw augmented data, $Y_{ij}^*$, as follows:

$$
Y_{ij}^* \quad
\begin{cases}
= Y_{ij}, & Z_{ij} = 0; \\
\sim \text{Truncated } N(X_{ij}, 1/\omega_j) \text{ on } (-\infty, M_{ij}], & Z_{ij} = 1; \\
\sim N(X_{ij}, 1/\omega_j), & Z_{ij} = 2.
\end{cases}
$$

2. For the process variable $\boldsymbol{X}$, we sample each process variable $\boldsymbol{X}_j$, $j = 1 \ldots, J$, conditional on the other process variables $\{\boldsymbol{X}_k : k \neq j\}$ and the other parameters in the model. We have that

$$
\begin{aligned}
& [\boldsymbol{X}_j | \{\boldsymbol{X}_k : k \neq j\}, \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{Y}^*, \boldsymbol{\omega}] \\
& \propto \quad [\boldsymbol{Y}_j^* | \boldsymbol{X}_j, \boldsymbol{\omega}] \, [\boldsymbol{X}_j | \{\boldsymbol{X}_k : k \neq j\}, \boldsymbol{\beta}_j, \tau_j] \prod_{l \in ch_j} [\boldsymbol{X}_l | \{\boldsymbol{X}_k : k \neq l\}, \boldsymbol{\beta}_l, \tau_l],
\end{aligned}
$$

where $ch_j$ denotes the columns of the process variable $\boldsymbol{X}$ that depend on $\boldsymbol{X}_j$ in their definition of the conditional mean $\mu_{ij}$ (i.e., those $l$ for which $s_{lk}^X = j$ for some $k$ or $s_{lm}^{EX} = j$ for some $m$).

3

The above is expression is equal to

$$\prod_{i=1}^{I} \left( \left[ Y_{ij}^* | X_{ij}, \omega_j \right] \left[ X_{ij} | \{ X_{ik} : k \neq j \}, \boldsymbol{\beta}_j, \tau_j \right] \prod_{l \in ch_j} \left[ X_{il} | \{ X_{ik} : k \neq l \}, \boldsymbol{\beta}_l, \tau_l \right] \right),$$

from which we can see that we can sample each $X_{ij}$, $i = 1, \ldots, I$, independently. Now

$$\left[ Y_{ij}^* | X_{ij}, \omega_j \right] \quad \propto \quad \exp\left( \frac{\omega_j}{2} (Y_{ij}^* - X_{ij})^2 \right),$$

$$\left[ X_{ij} | \{ X_{ik} : k \neq j \}, \boldsymbol{\beta}_j, \tau_j \right] \quad \propto \quad \exp\left( \frac{\tau_j}{2} (X_{ij} - \mu_{ij})^2 \right),$$

and for each $l \in ch_j$

$$\left[ X_{il} | \{ X_{ik} : k \neq l \}, \boldsymbol{\beta}_l \right] \propto \exp\left( \frac{\tau_l}{2} (X_{il} - \delta_{il}(j) - \epsilon_{il}(j))^2 \right). \tag{1}$$

In (1) we define

$$\delta_{ij}(o) \quad = \quad \alpha_j + \sum_{\{k : s_{jk}^X \neq o\}} \beta_{jk}^X X_{is_{jk}^X} + \sum_l \beta_{jl}^P P_{is_{jl}^P} + \sum_m \beta_{jm}^E E_{is_{jm}^E} + \sum_{\{n : s_{jn}^{EX} \neq o\}} \beta_{jn}^{EX} E_{is_{jn}^E} X_{is_{jn}^{EX}}$$

to be the conditional mean of $X_{ij}$ minus the regression terms that depend on $X_{io}$, and

$$\epsilon_{ij}(o) \quad = \quad \sum_{\{k : s_{jk}^X = o\}} \beta_{jk}^X X_{is_{jk}^X} + \sum_{\{m : s_{jm}^{EX} = o\}} \beta_{jm}^{EX} E_{is_{jm}^E} X_{is_{jm}^{EX}}$$

$$= \quad X_{io} \left[ \sum_{\{k : s_{jk}^X = o\}} \beta_{jk}^X + \sum_{\{m : s_{jm}^{EX} = o\}} \beta_{jm}^{EX} E_{is_{jm}^E} \right],$$

to be the regression terms that depend on $X_{io}$. Always, $\delta_{ij}(o) + \epsilon_{ij}(o) = \mu_{ij}$.

Rearranging the terms of $X_{ij}$, the distribution of $X_{ij}$ conditional on all other terms is $N(q/p, 1/p)$ where

$$q \quad = \quad \omega_j Y_{ij}^* + \tau_j \mu_{ij} + \sum_{l \in ch_j} \tau_l \left( X_{il} - \delta_{il}(j) \right) \left( \sum_{\{k : s_{jk}^X = n\}} \beta_{jk}^X + \sum_{\{m : s_{jm}^{EX} = n\}} \beta_{jm}^{EX} E_{is_{jm}^E} \right)$$

and

$$p \quad = \quad \omega_j + \tau_j + \sum_{l \in ch_j} \tau_l \left( \sum_{\{k : s_{jk}^X = n\}} \beta_{jk}^X + \sum_{\{m : s_{jm}^{EX} = n\}} \beta_{jm}^{EX} E_{is_{jm}^E} \right)^2.$$

4

3. For the regression parameters $\boldsymbol{\beta}$, we sample each $\boldsymbol{\beta}_j$ conditional on the other parameters in the model. Letting

$$
\boldsymbol{A}_j \;=\; \begin{bmatrix} 1 & \left\{X_{1s_{jk}^X}\right\} & \left\{P_{1s_{jk}^P}\right\} & \left\{E_{1s_{jk}^E}\right\} & \left\{E_{1s_{jk}^E}X_{1s_{jk}^{EX}}\right\} \\ 1 & \left\{X_{2s_{jk}^X}\right\} & \left\{P_{2s_{jk}^P}\right\} & \left\{E_{2s_{jk}^E}\right\} & \left\{E_{2s_{jk}^E}X_{2s_{jk}^{EX}}\right\} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & \left\{X_{Is_{jk}^X}\right\} & \left\{P_{Is_{jk}^P}\right\} & \left\{E_{Is_{jk}^E}\right\} & \left\{E_{Is_{jk}^E}X_{Is_{jk}^{EX}}\right\} \end{bmatrix},
$$

the distribution of $\boldsymbol{\beta}_j$ conditional on the other parameters is

$$
\mathrm{MVN}_{r_j}\left(\boldsymbol{R}_j^{-1}\boldsymbol{q}_j,\, \boldsymbol{R}_j^{-1}\right),
$$

where $\boldsymbol{R}_j = \tau_j \boldsymbol{A}_j^T \boldsymbol{A}_j + \boldsymbol{\Sigma}_j^{-1}$ and $\boldsymbol{q}_j = \tau_j \boldsymbol{A}_j^T \boldsymbol{X}_j + \boldsymbol{\Sigma}_j^{-1}\boldsymbol{m}_j$.

4. For each $j$, the conditional distribution of the process precision $\tau_j$, given the other parameters, is

$$
Ga\left(c_j + I/2,\; d_j + \sum_{i=1}^{N}(X_{ij} - \mu_{ij})^2/2\right),
$$

where $\mu_{ij}$ was defined above in Equation (1) of the paper.

# Models fit in "Effect and Pathways Modifiers in a Bayesian Pathways Analysis of the NHEXAS for Arsenic in EPA Region 5"

| UR | Variable | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | Clayton | Subpops | Reduced |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | **Intercept** | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ |
| | **PA** | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ |
| | **SD** | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ |
| | **SL** | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ |
| | **FD** | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ |
| | **BV** | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ |
| | **logcrtst** | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | | | | | | | | | |
| | mlogcrtn | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | ■ |
| | logcreatinine | | | | | | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | | | | | | | | | | ■ |
| | Female | ■ | | | | | | ■ | ■ | | | | | ■ | | ■ | ■ | ■ | ■ | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Male | | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | ■ |
| | Child4to11 | | | | | | | | | | | | | ■ | | | | ■ | | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age12to29 | | | | | | | | | | | | | ■ | ■ | ■ | ■ | ■ | ■ | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age30to39 | | | | | | | | | | ■ | | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age40to49 | | | | | | | | | | ■ | | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age50plus | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | tobacco | ■ | | ■ | | | | | | | | | | ■ | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | fruit | | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | |
| | hrs.home | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | ■ |
| | hrs.home**SD** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | | | | | | | | | | | | | |
| | hrs.home**SL** | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | ■ | | | | | | | | | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | |
| | fish | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | | | | | | | | | | |
| | new.fish | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | |
| | log hhold | | | | | | | | | | | | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | hhold | | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | ■ |
| | minushhold | | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | ■ | | | | | | |
| | hhold**SL** | | | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | |
| | log hhold**SL** | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Male**PA** | | | | | | | | | | | | | | | | | | | | | | | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | ■ |
| | Female**PA** | ■ | | | | | | ■ | ■ | | | | | ■ | | ■ | ■ | ■ | ■ | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age12to29**Female** | | | | | | | | | | | | | ■ | ■ | ■ | | ■ | | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age30to39**Female** | | | | | | | | | | ■ | | ■ | ■ | | | | ■ | | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age40to49**Female** | | | | | | | | | | ■ | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age50plus**Female** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Child4to11**PA** | | | | | | | | | | | | | ■ | | | | | | | | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Child4to11**BV** | | | | | | | | | | | | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Child4to11**FD** | | | | | | | | | | | | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Child4to11**SD** | | | | | | | | | | | | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Child4to11**SL** | | | | | | | | | | | | | ■ | | | | | | | | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | logBMI | | | | | | | | | | | ■ | | | | | | | | | | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age1 | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age2 | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age20to50 | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | age | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | race3_2 | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | race3_3 | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | race3_2**Age1** | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | race3_3**Age1** | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | race3_2**Age2** | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | race3_3**Age2** | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | race3_2**Age20to50** | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | race3_3**Age20to50** | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | age**Age20to50** | | | | | | | | ■ | ■ | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age12to19 | | | | | | | | | | ■ | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age20to29 | | | | | | | | | | ■ | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age50to59 | | | | | | | | | | ■ | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age60plus | | | | | | | | | | ■ | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age12to19**Female** | | | | | | | | | | ■ | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age20to29**Female** | | | | | | | | | | ■ | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age50to59**Female** | | | | | | | | | | ■ | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | Age60plus**Female** | | | | | | | | | | ■ | | ■ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

# Models fit in "Effect and Pathways Modifiers in a Bayesian Pathways Analysis of the NHEXAS for Arsenic in EPA Region 5"

Columns: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 | Clayton | Subpops | Reduced

**FD**
- Intercept
- IA
- SD
- WF
- tapwater
- tap.c
- no.fish
- fish
- new.fish
- fruit
- tap.c*WF
- tapwater*WF

**BV**
- Intercept
- IA
- SD
- WF
- tapwater
- tap.d
- glass
- tapwater*WF
- tap.d*WF
- glass*WF
- glass*tapwater
- glass*tapwater*WF

**IA**
- Intercept
- OA
- SL
- hhold
- minushhold
- hold*OA
- minushold*OA
- hhold*SL
- minushhold*SL
- no.c.air
- c.air
- no.c.air*OA
- c.air*OA

**OA**
- Intercept

**PA**
- Intercept
- OA
- IA
- tobacco
- gasequip
- workshop
- work.cat
- work20
- hrs.home
- hrs.home*OA
- hrs.home*IA
- Male
- tobacco1
- tobacco2

# Models fit in "Effect and Pathways Modifiers in a Bayesian Pathways Analysis of the NHEXAS for Arsenic in EPA Region 5"

Columns: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 Clayton Subpops Reduced

**SL**
- Intercept

**SD**
- Intercept
- IA

**WF**
- Intercept
- no.wellwater
- wellwater

**log crtst**
- Intercept
- logBMI
- Age12to29
- Age30to39
- Age40to49
- Age50plus
- Age1
- Age2
- Age20to50
- age
- race3_2
- race3_3
- Female
- race3_2*Age1
- race3_3*Age1
- race3_2*Age2
- race3_3*Age2
- race3_2*Age20to50
- race3_3*Age20to50