

Synthesizing Statistical and Mathematical Models to Address Non-Positivity

Paul Zivich
Assistant Professor
Department of Epidemiology
University of North Carolina at Chapel Hill

May 14, 2026

Acknowledgments¹

Funding: USA NIH: K01AI177102

Conflicts of Interest: None

Acknowledgments: Stephen Cole, Jess Edwards, Bonnie Shook-Sa, Eric Lofgren

Disclaimer: All errors are my own

Publication: Zivich et al. *J Epidemiol Community Health* 2026;80:352-356



pzivich@unc.edu



pzivich



pausalz@bsky.social

¹Footnotes are for asides or references

Question: What is the mean systolic blood pressure (SBP) in mm Hg among children and adolescents aged 2-17 in the United States?

Parameter: $\mu := E[Y]$

- Y_i : SBP for unit i
- $E[\cdot]$: expected value function

Data: NHANES 2017-2018

Systematic Error: SBP is missing for 44% of children

- R_i : indicator of SBP being observed

Extent of missingness means that estimate is sensitive to choices about missing data²

A complete-case analysis assumes marginal exchangeability

$$E[Y] = E[Y \mid R = 1]$$

²While a relatively simple descriptive question, we can consider causal inference to be a missing data problem and thus link the methods used here for confounding

A More Plausible Assumption

SBP is missing completely at random conditional on X

- X_i : age in years

$$E[Y | X = x] = E[Y | X = x, R = 1] \text{ for all values of } x$$

Under this assumption, μ is identified as³

$$E[Y] = E[E(Y | X = x, R = 1)]$$

or a weighted average of the complete-case means within strata of X

³Estimate with direct maximum likelihood estimator, g-computation

Positivity: all unique patterns of X have a nonzero probability of having Y being measured

$$\Pr(R = 1 \mid X = x) > 0 \text{ for all } x \in \mathcal{X}$$

Every unique age in the population has a non-zero probability of having Y observed

- Statement about population
- Deterministic or structural positivity

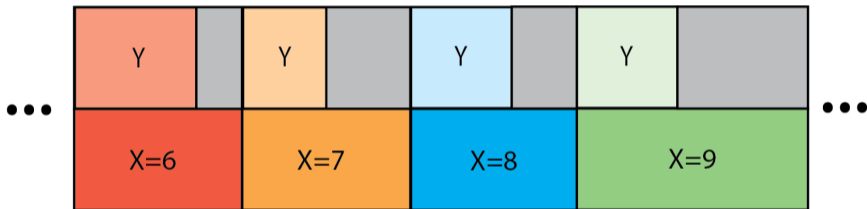
Formal: following definition of conditional expectation

$$E[Y \mid X = x, R = 1] = \frac{E[YR I(X = x)]}{\Pr(R = 1 \mid X = x) \Pr(X = x)}$$

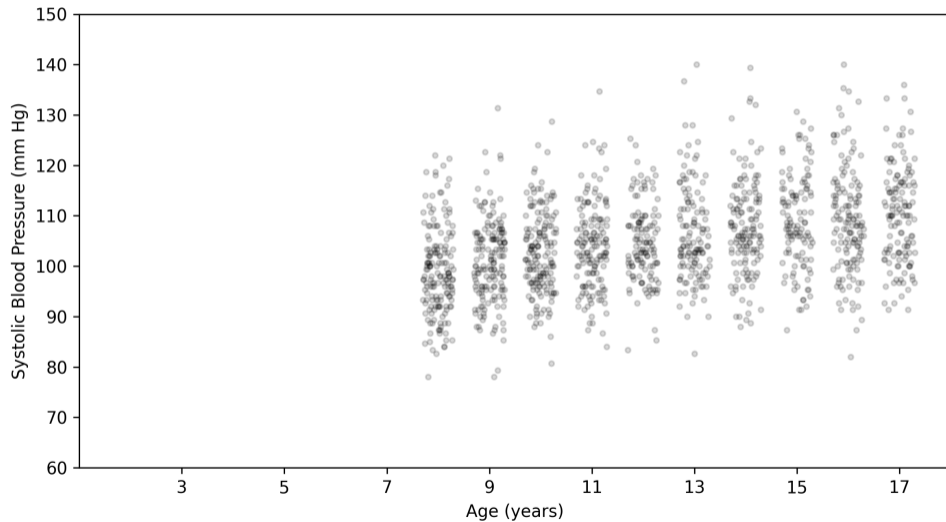
Intuitive: need some with $X = x$ to have Y observed at least sometimes so can 'stand-in' for those without Y measured

So, *positivity* comes along with *exchangeability*

Positivity for NHANES



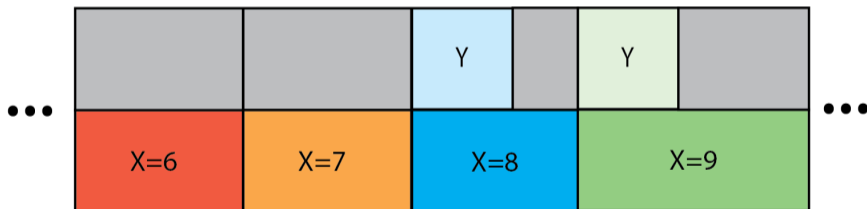
NHANES Data



Nonpositivity by Age in NHANES

In NHANES 2017-2018,

- “BP is measured on participants 8 years and older”⁴
- $\Pr(R = 1 \mid X = x) = 0$ for $x \in \{2, 3, 4, 5, 6, 7\}$



⁴https://www.cdc.gov/Nchs/Nhanes/2017-2018/BPX_J.htm

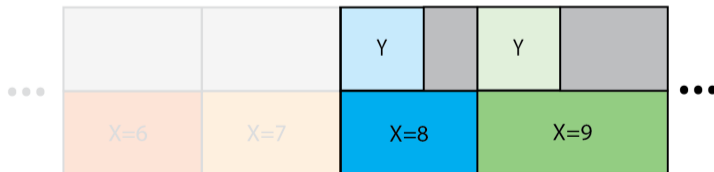
Addressing Positivity Violations

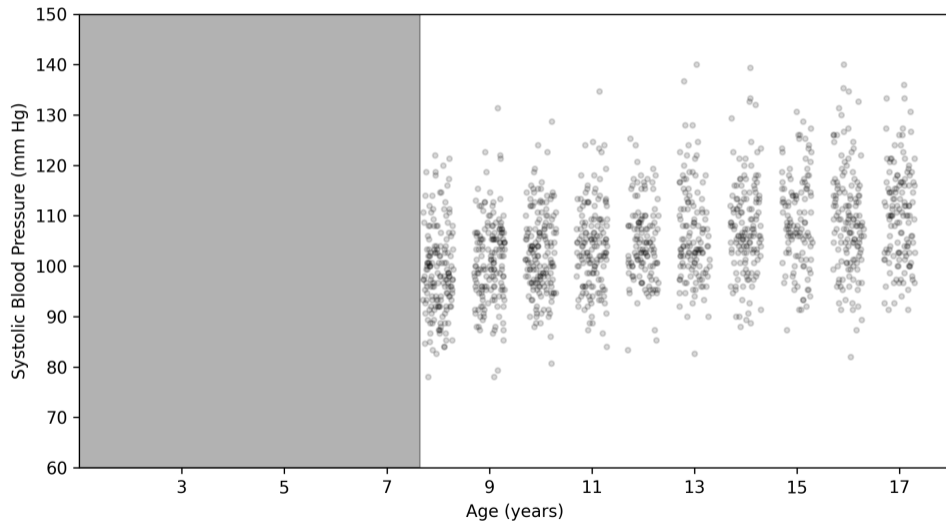
Option 1: Change Parameter

Parameter (revised): mean SBP for those 8 or older

$$E[Y \mid X \geq 8]$$

Question: what is the mean systolic blood pressure (SBP) in mm Hg among children and adolescents aged 8-17 in the US?



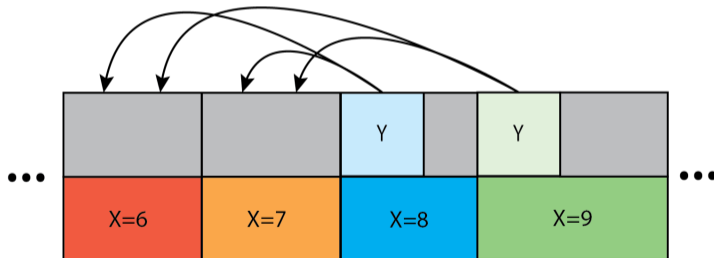


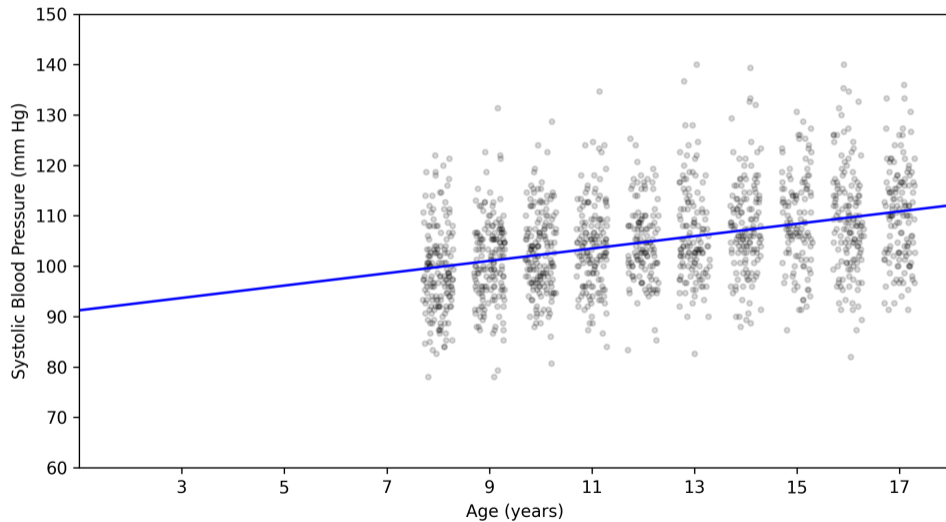
Option 2. Extrapolate from Statistical Model

Fit a parametric model

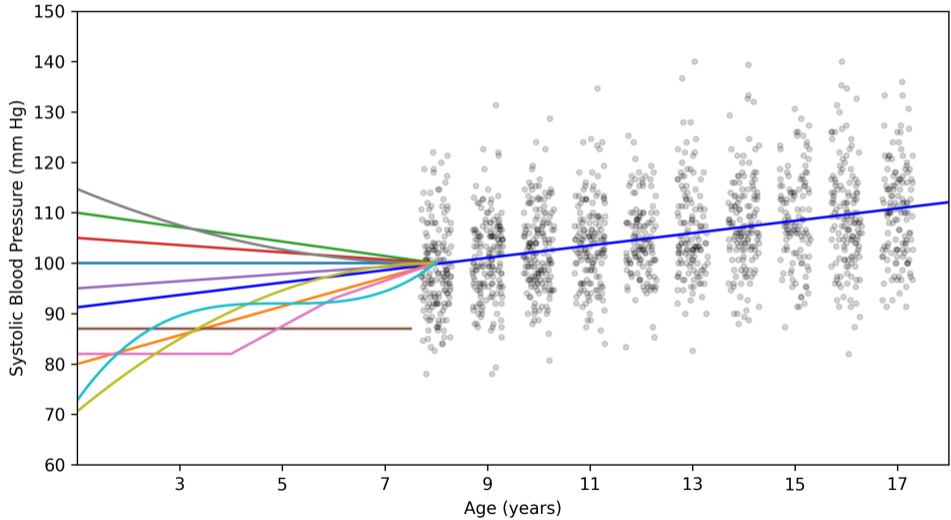
$$E[Y | X \geq 8, R = 1] = \beta_0 + \beta_1 X$$

then fill in missing SBP with model predictions





NHANES Data



2. Extrapolate

Use of this model has two assumptions:

- a. Linear relationship with age for those 8-17
- b. Linear relationship extends to those 2-7

(a) is in-principle checkable

(b) is not checkable under nonpositivity

Option 3: Bounds and Synthesis Modeling

Divide parameter of interest in parts

$$\mu = E[Y | X \geq 8] \Pr(X \geq 8) + E[Y | X < 8] \Pr(X < 8)$$

Overall mean

Mean in nonpositive region

Mean in positive region

Here,

$$E[Y | X \geq 8] = E[E(Y | X, R = 1, X \geq 8) | X \geq 8]$$

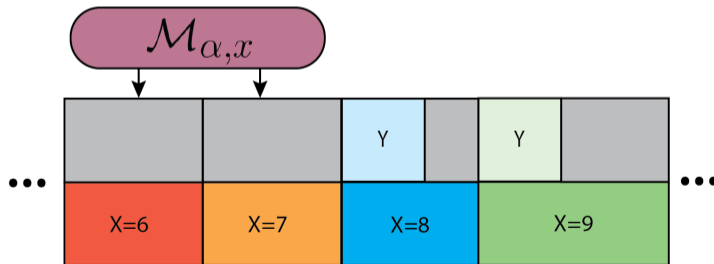
from conditional exchangeability in positive region

Option 3: Bounds and Synthesis Modeling

Available data leaves us stuck with $E[Y | X < 8] = ?$

To progress, will appeal to external information for

- Bounds
- Sensitivity Analysis
- Point Identification



Option 3a: Bounds

In our context, know that SBP must be a (finite) positive number, so

$$E[Y | X = x] \in (0, \infty) \text{ for all } x \in \{2, 3, 4, 5, 6, 7\}$$

Context allows us to further narrow the plausible region to

$$E[Y | X = x] \in [140, 200] \text{ for all } x \in \{2, 3, 4, 5, 6, 7\}$$

as these boundary values would constitute medical emergencies

Therefore,

$$50 + \Pr(X \geq 8) \left\{ E[Y | X \geq 8] - 50 \right\} \leq \mu \leq 200 + \Pr(X \geq 8) \left\{ E[Y | X \geq 8] - 200 \right\}$$

Option 3b: Sensitivity Analysis

Rather than evaluate at the extremes, evaluate

$$\mu = E[E(Y | X, R = 1, X \geq 8) | X \geq 8] \Pr(X \geq 8) + q \Pr(X < 8)$$

where $q \in [50, 200]$

Assess how μ varies as a function of the possible means for the nonpositive region

- May help to further contextualize the bounds

Option 3c: Leverage a Mathematical Model

Previous only provide range of equally plausible values

- Can leverage more specific information when available
- Published distributions of SBP by age, gender, height⁵

Use this additional information to construct a mathematical model that produces a probability distribution $\mathcal{M}_{\alpha,x}$ for age x given parameters α , and assume

$$E[Y | X = x] \in \mathcal{M}_{\alpha,x} \text{ for all } x \in \{2, 3, 4, 5, 6, 7\}$$

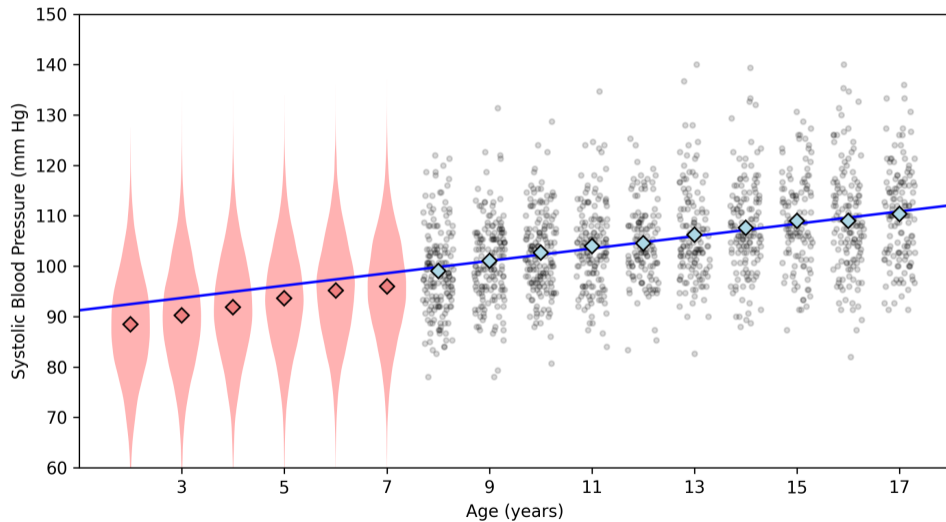
⁵Flynn et al. 2017 *Pediatrics* 140(3):e20171904

Then can combine statistical and mathematical models for identification

$$\begin{aligned} \mu &= E[E(Y | X, R = 1, X \geq 8) | X \geq 8] \Pr(X \geq 8) \\ &+ E[E(\mathcal{M}_{\alpha,x} | X, X < 8) | X < 8] \Pr(X < 8) \end{aligned}$$

Build $\mathcal{M}_{\alpha,x}$ from published age-, gender-, and height-specific SBP distributions

NHANES Data



Option 3a/b: estimate $E[E(Y | X, R = 1, X \geq 8) | X \geq 8]$ using statistical modeling,⁶ set q , solve for μ .⁷

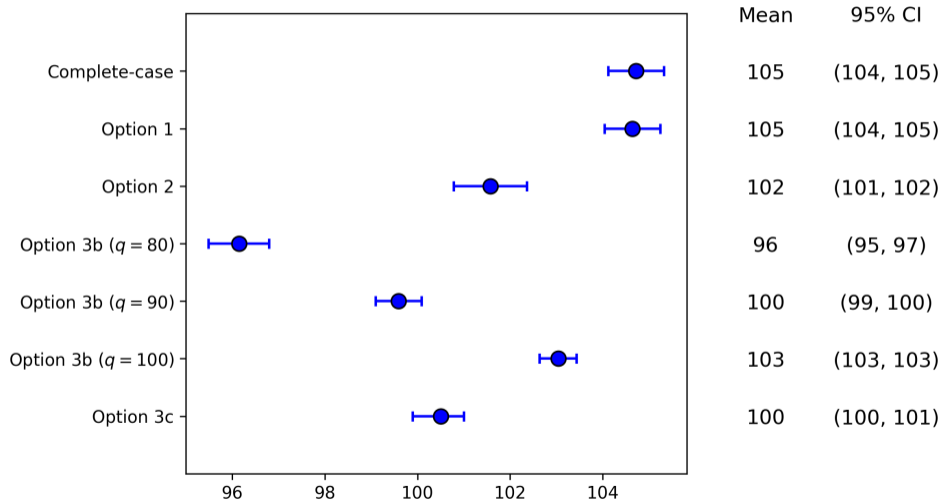
Option 3c: Monte Carlo procedure to capture uncertainty

1. Resample with replacement observed data
2. Estimate a statistical model for $E(Y | X, R = 1, X \geq 8)$ and then impute Y for observations with $X \geq 8$
3. For those $X < 8$, impute Y using a random draw from the mathematical model covariate-specific distributions
4. Take the mean of the imputed Y for all observations

⁶Here, I use g-computation for ease but could also use IPW or AIPW. The publication applies all three

⁷Use the pointwise Confidence Intervals which are conservative, Richardson et al. *Stat Sci* 2014;29(4):596-618

Results by Method



Conclusions

Positivity is often neglected as an assumption

- But it remains vital for identification of many parameters

Restricting parameters to positive regions often reasonable

- But may limit utility of research
- Extrapolation premised on strong assumptions
- Synthesis modeling allows for exploration and weakening of extrapolation assumptions

Compare statistical and mathematical models for positive region

- May provide some assurance for mathematical model
- Regularize statistical model for efficiency

Exchangeability (and thus positivity) is used to address other systematic errors

- Approach can be used to address nonpositivity in these settings⁸

Analysis of mathematical models and improving their robustness⁹

⁸Originally proposed in the context of transporting trials: Zivich et al. *Epidemiology* 2024;35(1):23-31, Zivich et al. *JRSSA* 2025;188(1):158-180

⁹Some work on formalizing identification with mathematical models: Zivich *arXiv:2511.01960*

Publication: Zivich et al. *J Epidemiol Community Health* 2026;80:352-356



pzivich@unc.edu



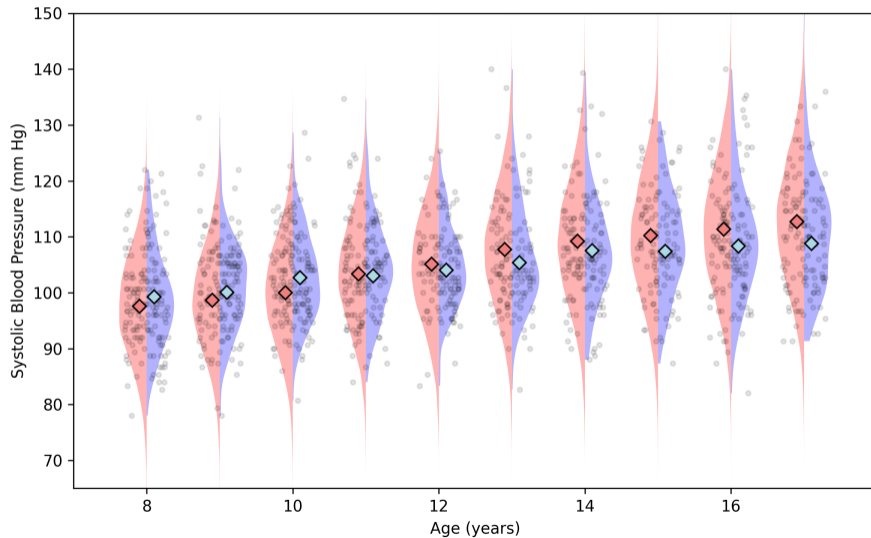
pzivich



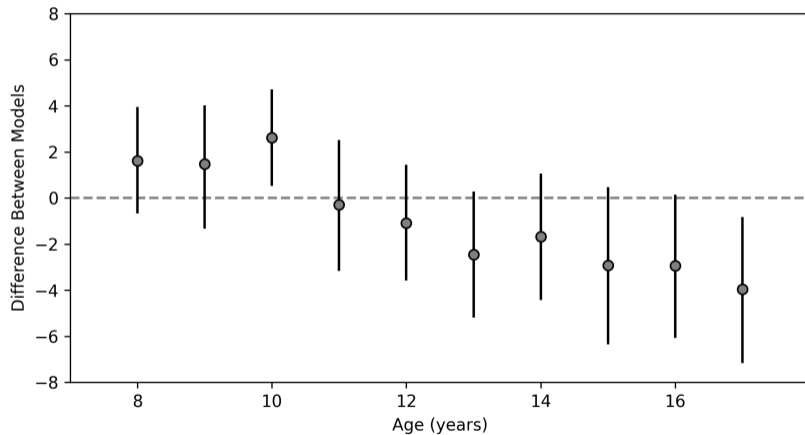
pausalz@bsky.social

Appendix

Comparing Statistical and Mathematical Models



Comparing Statistical and Mathematical Models¹⁰



¹⁰Intervals are pointwise 95% confidence intervals, so anti-conservative for overall comparison