

A BNP Model for Zero-Inflated Outcomes with Applications in Causal Inference

Arman Organisian, Jason Roy, Nandita Mitra

Department of Biostatistics, Epidemiology, and Informatics
Division of Biostatistics
University of Pennsylvania

JSM 2018
July 30, 2018

BASIC SETTING

Consider a cross-sectional study with

- ▶ Binary treatment: $A \in \{0, 1\}$
- ▶ Continuous outcome: $Y \in (-\infty, \infty)$
- ▶ Single, continuous confounder: $L \in (-\infty, \infty)$
- ▶ Goal: estimate Ψ - the average causal effect of A on Y

$$\Psi = E[Y^{A=1} - Y^{A=0}]$$

If standard causal assumptions (ignorability, consistency, positivity, no interference) are met, can use Standardization (Robins, 1986)

$$E[Y^{A=a}] = \int E[Y|A=a, L; \beta] dF(L; \alpha)$$

BAYESIAN STANDARDIZATION

Terms of Ψ are computed using the posterior predictive distribution (Keil, 2017),

$$E[\tilde{y}^a | Y, L] = \int_{\alpha} \int_{\beta} \int_{\tilde{L}} E[\tilde{y} | A = a, \tilde{L}, \beta] p(\tilde{L} | \alpha) p(\beta, \alpha | Y, L) d\tilde{L} d\beta d\alpha \quad (1)$$

Need to model conditional distribution of Y and distribution of L . E.g.,

$$E[Y | A = a, L = l] = \beta_0 + \beta_1 a + \beta_2 l$$

- ▶ Imputation model $E[Y | A = a, L, \beta]$ needs to be correctly specified.
- ▶ Two sets of rigid assumptions: causal assumptions and statistical assumptions.
- ▶ Flexible nonparameteric methods can at least help us relax the latter.
 - ▶ Especially important for modeling cost data.

DP MIXTURE OF ZERO-INFLATED REGRESSION

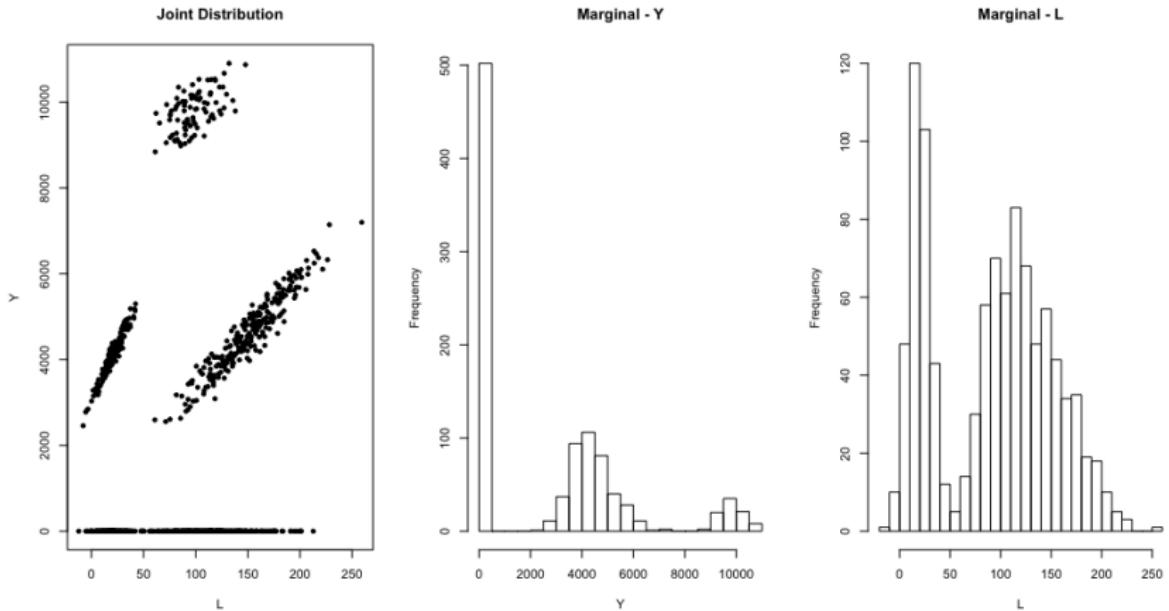
Building off of previous methods (Hannah, 2011) (Roy, 2018), we propose the following generative model

$$\begin{aligned} y_i | z_i, \mathbf{x}_i &\sim \begin{cases} \delta_0(y_i), & z_i = 1 \\ N(\mathbf{x}_i^T \boldsymbol{\beta}_i, \phi_i), & z_i = 0 \end{cases} \\ z_i | \mathbf{x}_i &\sim Ber(\expit(\mathbf{x}_i^T \boldsymbol{\gamma}_i)) \\ x_{i,j} &\sim g_{j,i} = \begin{cases} N(\lambda_{j,i}, \tau_{j,i}), & x \text{ is continuous} \\ Ber(p_{j,i}), & x \text{ is binary} \end{cases}, \quad j \in \{1, 2, \dots, d\} \\ (\boldsymbol{\beta}_i, \boldsymbol{\gamma}_i, \boldsymbol{\lambda}_i, \boldsymbol{\tau}_i, \mathbf{p}_i) | G &\sim G \\ G &\sim DP(\alpha G_0) \end{aligned} \tag{2}$$

- ▶ Nonparametric in the sense that there are infinitely many potential parameters.
- ▶ But DP prior induces clustering. Infinitely many possible clusters.



SOME SIMULATED DATA



Center for
Causal Inference
 $C \rightarrow C \rightarrow I$



CCEB

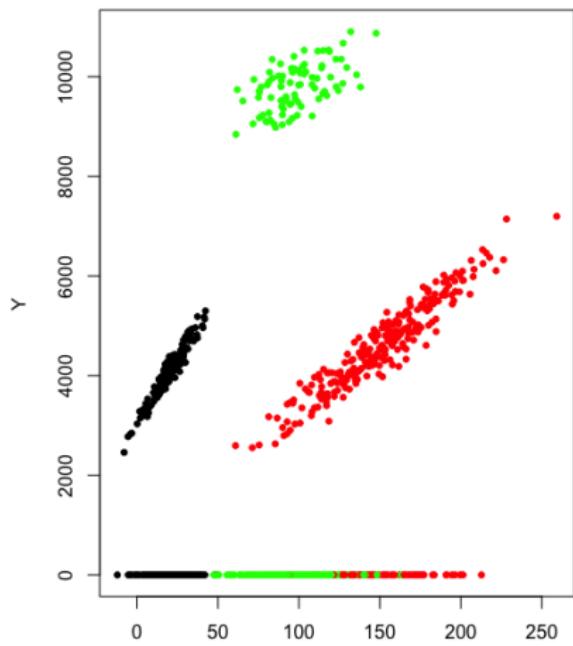
Perelman
School of Medicine
UNIVERSITY OF PENNSYLVANIA

ESTIMATION USING MCMC METHODS

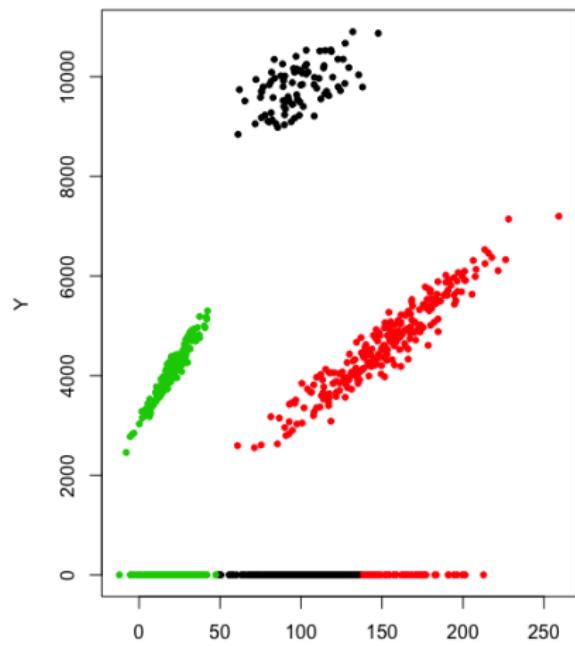
We use a Metropolis-within-Gibbs approach extended from (Neal, 2000) and similar to (Roy, 2018).

DP-INDUCED CLUSTERING

True Class Membership



Posterior Mode Class Membership



STANDARDIZATION - DRAWING FROM POSTERIOR PREDICTIVE UNDER DIFFERENT INTERVENTIONS

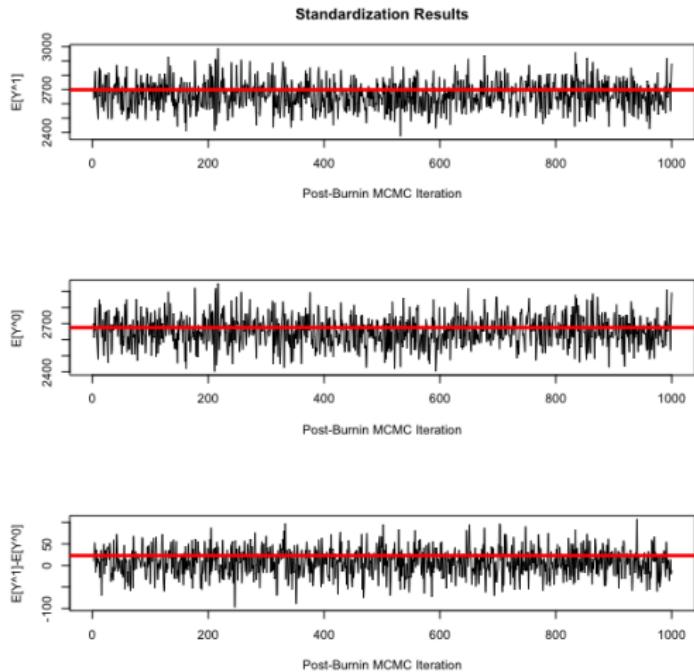
$$\begin{aligned} p(\tilde{y}^a | \mathbf{y}, \mathbf{x}, \mathbf{z}) &= \sum_{k=1}^{\infty} \int_{\boldsymbol{\theta}_{x,k}} \int_{\boldsymbol{\beta}_k} \int_{\phi_k} \int_{\boldsymbol{\gamma}} \int_{\tilde{\mathbf{x}}} \sum_{l \in \{0,1\}} p(\tilde{y} | \boldsymbol{\beta}_k, \phi_k, c = k, \tilde{\mathbf{z}} = l, \tilde{\mathbf{x}}^a) \cdot p(\tilde{\mathbf{z}} = l | c = k, \tilde{\mathbf{x}}^a, \boldsymbol{\gamma}_k) \\ &\quad \cdot p(\tilde{\mathbf{x}}^a | \boldsymbol{\theta}_{x,k}, c = k) \cdot p(\boldsymbol{\beta}_k, \phi_k, \boldsymbol{\gamma}_k, \boldsymbol{\theta}_{x,k}, c = k | \mathbf{y}, \mathbf{x}, z = l) d\tilde{\mathbf{x}} d\boldsymbol{\gamma}_k d\phi_k d\boldsymbol{\beta}_k d\boldsymbol{\theta}_{x,k} \end{aligned} \tag{3}$$

$$\Psi = E[p(\tilde{y}^1 | \mathbf{y}, \mathbf{x}, \mathbf{z})] - E[p(\tilde{y}^0 | \mathbf{y}, \mathbf{x}, \mathbf{z})]$$

- ▶ Developed Monte Carlo procedure for evaluation of this integral.
- ▶ Can compute other causal contrasts, e.g. $E[Y^1]/E[Y^0]$, easily.
- ▶ Can compute conditional causal effects using appropriately conditional posterior predictive.
- ▶ Interval estimates constructed in the usual ways.



STANDARDIZATION - MCMC CHAINS



REFERENCES

- ▶ Roy Jason, Lum Kirsten J., Zeldow Bret, Dworkin Jordan D., Re Vincent Lo, and Daniels Michael J. Bayesian nonparametric generative models for causal inference with missing at random covariates. *Biometrics*, 0(0).
- ▶ Alexander P Keil, Eric J Daza, Stephanie M Engel, Jessie P Buckley, and Jessie K Edwards. A bayesian approach to the g-formula. *Statistical Methods in Medical Research*, 0(0):0962280217694665, 0. PMID: 29298607.
- ▶ Lauren A Hannah, David M Blei, and Warren B Powell. Dirichlet process mixtures of generalized linear models. *Journal of Machine Learning Research*, 12(Jun):1923?1953, 2011.
- ▶ Radford M. Neal. Markov chain sampling methods for dirichlet process mixture models. *Journal of Computational and Graphical Statistics*, 9(2):249?265, 2000.

APPENDIX 1: CAUSAL ASSUMPTIONS

- ▶ Ignorability: $Y_i^{A_i=a} \perp A_i = a | L_i$. Conditional on observed covariates, potential cost is independent of treatment assignment. In randomized control trials, this conditional independence holds by virtue of randomization.
- ▶ Consistency: cost Y_i observed under the actual treatment $A_i = a$ is equal to $Y_i^{A_i=a}$. Specifically, $Y_i^{A_i=a} = Y_i | A_i = a$.
- ▶ No interference: one patient's treatment assignment does not impact another's potential outcome - $Y_i^{A_i=a} \perp A_j, \forall i \neq j$. A common example of interference is a setting in which Y represents someone's infection status and A represents someone's vaccination status against the disease in question. It is reasonable to consider that one patient's vaccination status may impact another's infection status.
- ▶ Positivity: no patient has a deterministic treatment. That is, the probability is strictly between 0 and 1 $0 < P(A_i = 1 | L_i) < 1$. If this assumption did not hold, then one of subject i 's potential outcomes would be undefined.