

# Bayesian Nonparametric Model for Zero-Inflated Outcomes

## Clustering, Prediction, and Causal Inference

Arman Oganisian

with Jason Roy and Nandita Mitra

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

July 31, 2019

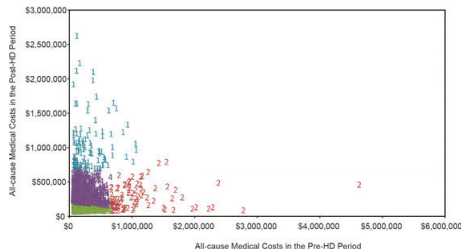
1/16



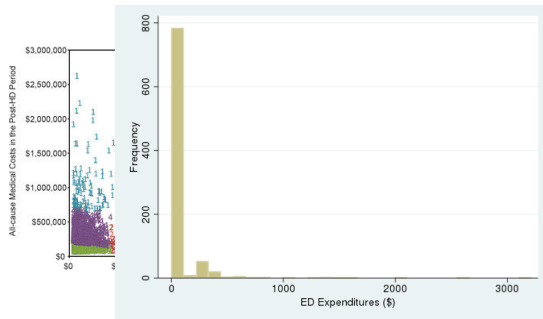
# MOTIVATION

- ▶ Health policy questions involving costs are complicated:
  - ▶ **Causality**: How different would average costs have been under alternative treatment?
  - ▶ **Prediction**: How much medical costs will subject  $X$  likely accumulate?
  - ▶ **Clustering**: Can we identify interesting patient subgroups?
- ▶ Cost data are complicated:
  - ▶ zero-inflation
  - ▶ skewness
  - ▶ multimodality
- ▶ Complicated questions with complicated data.

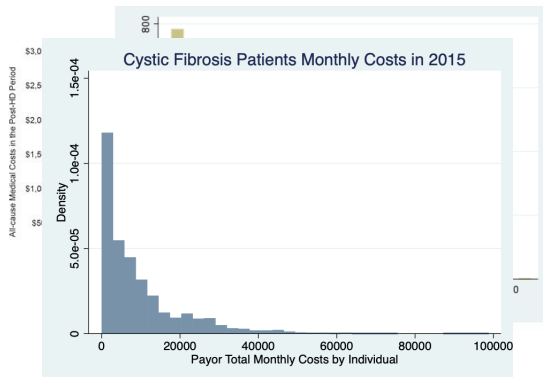
# EXAMPLES FROM LITERATURE



# EXAMPLES FROM LITERATURE



# EXAMPLES FROM LITERATURE



# WHAT ARE OUR OPTIONS?

Observed data  $D = (D_i)_{i=1:n} = (Y_i, x_i = (A_i, L_i))_{i=1:n}$

- ▶ Ignore zeros (efficiency loss, underestimate treatment effect).
- ▶ Add a penny. (ad hoc, structural zeros  $\rightarrow$  structural pennies)
- ▶ Hurdle model (parametric, no clustering)

$$Y_i \mid A_i, L_i \sim \pi(x_i' \gamma) \delta_0(y_i) + (1 - \pi(x_i' \gamma)) \cdot f(y_i \mid x_i' \beta)$$

# BAYESIAN STANDARDIZATION

Under certain identification assumptions, can compute  $\Psi = E[Y^1 - Y^0]$ .

$$E[Y^a|D] = \int_{\beta} \int_L E[Y|A = a, L, \beta] p(L) p(\beta|D) dL d\beta$$

# ZERO-INFLATED DIRICHLET PROCESS MIXTURE

Generative model for the full joint data  $p(D_i|\omega_i)$

$$Y_i \mid A_i, L_i \sim \pi (x_i' \gamma_i) \delta_0 (y_i) + (1 - \pi (x_i' \gamma_i)) \cdot N (y_i \mid x_i' \beta_i, \phi_i)$$

$$A_i \mid L_i \sim \text{Ber} (\text{expit} (m_i' \eta_i))$$

$$L_i \sim p (l_i \mid \theta_i)$$

Joint prior for  $\omega_i = (\beta_i, \phi_i, \gamma_i, \eta_i, \theta_i)$

$$\omega_i \mid G \sim G$$

$$G \sim DP (\alpha G_0)$$



# INDUCED POSTERIOR CLUSTERING

Conditional posterior of  $i^{th}$  subject's parameters

$$p(\omega_i | \omega_{1:(i-1)}, D) \propto \frac{\alpha}{\alpha + i - 1} p(D_i | \omega_i) G_0(\omega_i) + \frac{1}{\alpha + i - 1} \sum_{j < i} p(D_i | \omega_j) I(\omega_i = \omega_j)$$

- ▶ Data adaptive.
- ▶ Posterior clustering.
- ▶ Flexible predictions by ensembling cluster-specific models.

# MCMC INFERENCE: PARTITIONS AND PREDICTIONS

8/16



Center for  
Causal Inference @StableMarkets 



DEPARTMENT of  
BIostatISTICS  
EPIDEMIOLOGY &  
INFORMATICS

 Perelman  
School of Medicine  
UNIVERSITY of PENNSYLVANIA

# SOME RESULTS

In the paper we have

- ▶ expressions for key posterior distributions.
- ▶ required Monte Carlo procedures.
- ▶ standardization procedure around the model.
- ▶ posterior predictive checks assessing positivity.
- ▶ hard posterior classification in presence of label switching.
- ▶ uncertainty visualization in cluster assignment.

# SIMULATIONS RESULTS

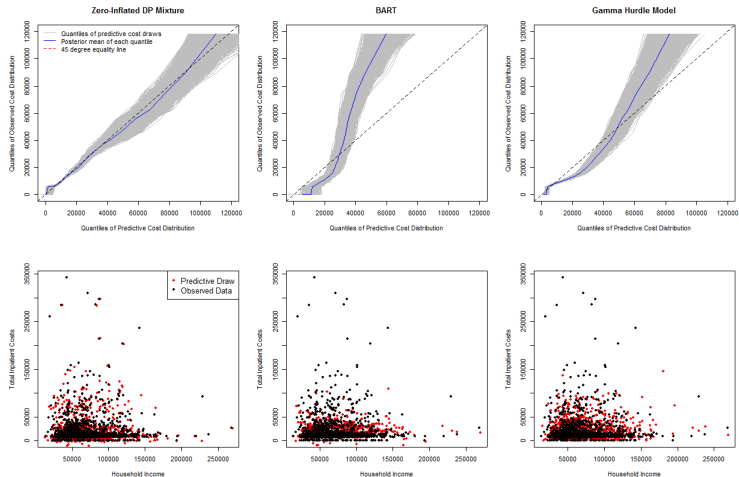
- ▶ Simulated cost data from three distinct Bernoulli-Gamma hurdle model.
- ▶  $n = 3,000$  subjects, 5 covariates, single binary treatment.

| DGP        | Model            | Bias  | Coverage | Rel. Interval Width |
|------------|------------------|-------|----------|---------------------|
| Clustered  | Zero-Inflated DP | -.081 | 94.3%    | 1.10                |
|            | BART             | -.746 | 76.2%    | 1.34                |
|            | Doubly Robust    | .795  | 87.1%    | 1.70                |
|            | Gamma Hurdle     | -.509 | 79.8%    | 1                   |
|            | Gamma +.01       | 1.817 | 4.7%     | 1.39                |
| Parametric | Zero-Inflated DP | .097  | 95.1%    | 1.01                |
|            | BART             | -.054 | 96.1%    | 1.09                |
|            | Doubly Robust    | -.027 | 95.9%    | 1.07                |
|            | Gamma Hurdle     | -.014 | 95.1%    | 1                   |
|            | Gamma +.01       | -.489 | 100%     | 2.32                |

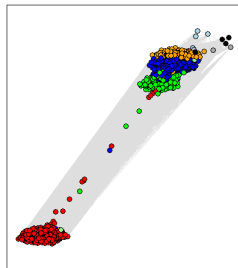
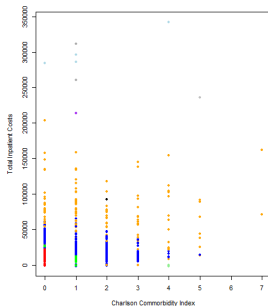
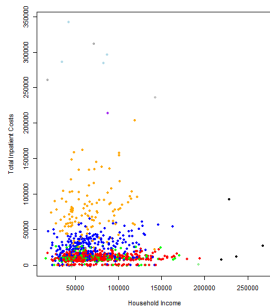
# TREATMENT COSTS FOR ENDOMETRIAL CANCER

- ▶ Data source: SEER-Medicare.
- ▶ Endometrial cancer patients ( $N \approx 1,000$ ).
- ▶ Treatment: post-hysterectomy **radiation** vs. **chemotherapy**.
- ▶ Outcome: Total inpatient costs over 2 years.
  - ▶ Skewed, zero-inflated
  - ▶ Chemo arm: 15% zeros; RT arm: 8%
- ▶ Measured confounders: tumor grade, cancer stage, CCI.

# PREDICTIONS



# POSTERIOR CLUSTERING

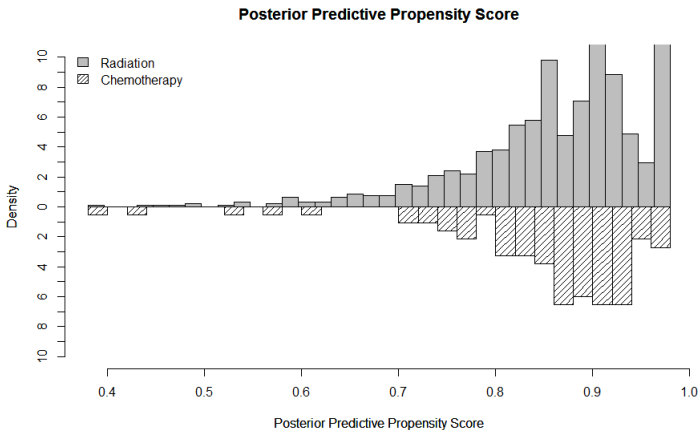


# AVERAGE TREATMENT EFFECTS

|                  | Avg. Causal Effect              | Median Causal Effect         | Zero - Risk Ratio     |
|------------------|---------------------------------|------------------------------|-----------------------|
| Zero-Inflated DP | 1672.62<br>(-2566.42, 5722.56)  | 872.68<br>(-833.35, 2790.18) | 0.498<br>(0.31, 0.78) |
| BART             | 1779.62<br>(-6085.89, 9797.13 ) | -                            | -                     |
| Gamma Hurdle     | 2016.71<br>(-1499.38, 5593.40)  | -                            | .505<br>(.34, .76)    |
| Gamma +.01       | 4889.00<br>(1004.37, 8795.61)   | -                            | -                     |



# POSTERIOR PREDICTIVE PROPENSITY SCORES



# PAPER, TUTORIAL, SOFTWARE

- ▶ arXiv: <https://arxiv.org/abs/1810.09494>
- ▶ Interactive DP Tutorial with R Shiny: <https://stablemarkets.shinyapps.io/dpmixapp/>
- ▶ ChiRP R package: <https://stablemarkets.github.io/ChiRPsite/index.html>

arXiv



Tutorial



ChiRP

