

# Nonparametric Bayes and Causal Sensitivity Analysis

A.W. van der Vaart

TU Delft

Lund Statistics Day 2024

joint with



Bart Eggen, TU Delft



Stephanie van der Pas, AMC



Chris van Vliet, TU Delft

# Causal Inference

$A \in \{0,1\}$  treatment indicator

$y^1, y^0$

outcomes if  $A=1,0$ . "counterfactual"

Causal effect  $E y^1 - E y^0$

# Causal Inference

$A \in \{0,1\}$  treatment indicator

$y^1, y^0$  outcomes if  $A=1,0$ . "counterfactual"

Causal effect  $E y^1 - E y^0$  only estimable under conditions

$Z$  covariates

Observed  $(A, y, z)$ ,  $y = y^A$ .

LEM if  $y^1, y^0 \perp\!\!\!\perp A \mid Z$ , then

$$E y^1 - E y^0 = E_Z (E(y \mid A=1, Z) - E(y \mid A=0, Z))$$

LEM If  $y^1, y^0 \perp\!\!\!\perp A \mid Z$ , then

$$E y^1 - E y^0 = E_Z \left( E(y^1 \mid A=1, Z) - E(y^1 \mid A=0, Z) \right)$$

proof

$$\begin{aligned} E y^a &= E_Z E(y^a \mid Z) \\ &= E_Z E(y^a \mid A=a, Z) && \text{(assumption)} \\ &= E_Z E(y \mid A=a, Z) && \text{(consistency)} \end{aligned}$$

(need also that conditioning on  $\{A=a, Z\}$  makes sense:

positivity:  $P(A=a \mid Z) > 0$ )

□



# Sensitivity Analysis

$$y^1, y^0 \perp\!\!\!\perp A \mid Z$$

"Conditional Exchangeability"  
or

"No unmeasured confounding"

$Z$  has to be rich enough to make this true  
(and not too rich)

Sensitivity Analysis: consider deviations

# Sensitivity Analysis

$$y^1, y^0 \perp\!\!\!\perp A \mid Z$$

"Conditional Exchangeability"  
or

"No unmeasured confounding"

$Z$  has to be rich enough to make this true  
(and not too rich)

Sensitivity Analysis: consider deviations

## Two approaches

- Assume  $\exists u$  with  $y^1, y^0 \perp\!\!\!\perp A \mid Z, u$
- Model  $A \mid Z, y^a$ ,  $a \in \{0, 1\}$

# Bayesian Sensitivity Analysis

Both approaches add an unidentifiable  
"sensitivity parameter"

Causal effect is identifiable only given the sensitivity parameter.

Bayesian approach puts a prior on this parameter and obtains a posterior of the causal effect

As usual, except

- posterior does not contract to point and
- priors matter

Two approaches to model failure of  $y^1, y^0 \perp\!\!\!\perp A \mid Z$

• Assume  $\exists U$  with  $y^1, y^0 \perp\!\!\!\perp A \mid Z, U$

• Model  $A \mid Z, y^a$ ,  $a \in \{0, 1\}$

$$\text{logit } P(Y=1 | A, Z, U) = \beta_0 + \beta_1 A + \alpha U + \eta^T Z$$

$$\text{logit } P(U=1 | A, Z) = \delta_0 + \delta_1 A + \xi^T Z$$

$A$  beta blocker  
 $Y$  1-year mortality  
 $U$  unmeasured  
 $Z$  covariates

$$Y^1, Y^0 \perp\!\!\!\perp A \mid Z, U \Rightarrow \text{logit } P(Y^a=1 | A=a, Z, U) = \beta_0 + \beta_1 a + \alpha U + \eta^T Z$$

"causal effect"  $\leftrightarrow \beta_1$

get posterior of  $\beta_1$  by Gibbs sampling with  $U$  missing:

repeat  $\beta_0, \beta_1, \alpha, \eta, \delta_0, \delta_1 \mid \text{DATA}, U$

$U \mid \text{DATA}, \beta_0, \beta_1, \alpha, \eta, \delta_0, \delta_1$

$$y | A, Z, U \sim N(f(A, Z) + \alpha_1 U, \sigma^2)$$

$$P(A=1 | Z, U) = \Phi(Z\beta + \alpha_2 U)$$

$$U | Z, U \sim \text{Bernoulli}(\pi)$$

prior on  $f$  : BART

$$y^1, y^0 \perp\!\!\!\perp A | Z, U \Rightarrow E y^a = E E(y | A=a, Z, U) = E f(a, Z) + \alpha_1 E U$$

$$E y^1 - E y^0 = E f(1, Z) - E f(0, Z).$$

get posterior of  $E y^1 - E y^0$  from posterior of  $f$  by  
Gibbs sampling with  $U$  missing

# Nonparametric Bayes (1): BART

Chipman, George,  
McCulloch, 1998.

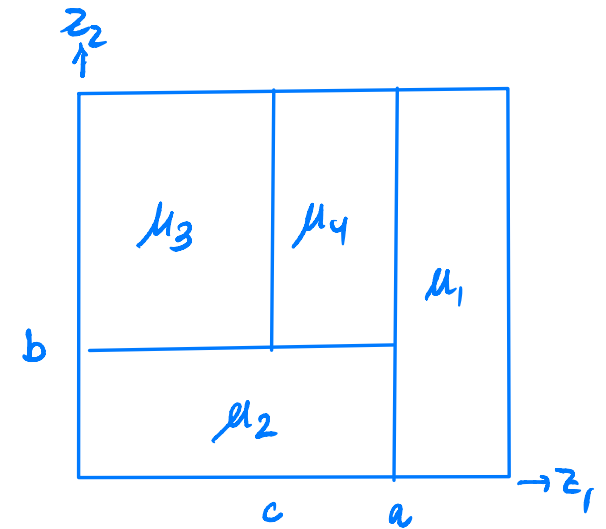
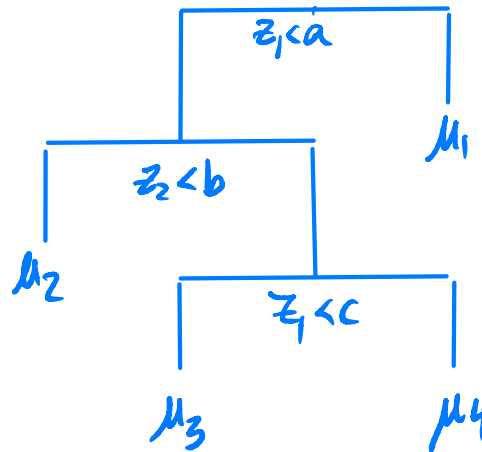
Bayesian Additive Regression Tree

$$f(z) = g(z; T_1, M) + g(z; T_2, M) + \dots + g(z; T_m, M_m)$$

$$z \mapsto g(z; T, M)$$

$T = \text{tree}$

$$M = (a, b, c, \dots, \mu_1, \mu_2, \dots)$$



## PRIOR

- Node at level  $h$  is split with probability  $0.95/(1+h)^2$
- Splitting variable  $z_i$  and value  $a, b, c, \dots$  chosen uniformly.
- Leaf values  $\mu_1, \mu_2, \dots \stackrel{\text{iid}}{\sim} N(0, \sigma^2/m)$ .

+ empirical Bayes

Two approaches to model failure of  $y^1, y^0 \perp A \mid Z$

• Assume  $\exists u$  with  $y^1, y^0 \perp A \mid Z, u$

• Model  $A \mid Z, y^a$ ,  $a \in \{0, 1\}$



Robins, Rotnitzky, Scharfstein, 2000  
Scharfstein, Daniels, Robins, 2003

---

FOCUS ON  $y^1$   
 $y^0$  ANALOGOUS

$$\text{logit } P(A=1 | y^1, z) = \eta(z) + \rho(y^1 | z)$$

sensitivity model

FOCUS ON  $y^1$   
 $y^0$  ANALOGOUS

$$\text{logit } P(A=1 | y^1, z) = \eta(z) + q(y^1 | z)$$
$$P(y^1 = 1 | z) = P(\cdot | z)$$

sensitivity model  
outcome

FOCUS ON  $y^1$   
 $y^0$  ANALOGOUS

$$\text{logit } P(A=1 | y^1, z) = \eta(z) + q(y^1 | z)$$

$$P(y^1 e \cdot | z) = P(\cdot | z)$$

sensitivity model

outcome

---

$$P(A=1 | z) = \pi(z)$$

$$P(y^1 e \cdot | A=1, z) = P_1(\cdot | z)$$

propensity score

observed outcome

---

$$\text{logit } P(A=1|y', z) = \eta(z) + q(y'|z)$$

$$P(y^1 e \cdot | z) = P(\cdot | z)$$

sensitivity model  
outcome

$$P(A=1|z) = \pi(z)$$

$$P(y^1 e \cdot | A=1, z) = P_1(\cdot | z)$$

propensity score  
observed outcome

$$P(y^1 e \cdot | A=0, z) = P_0(\cdot | z)$$

unobserved outcome

LEM •  $\forall q \exists$  bijections  $(\eta, P) \leftrightarrow (\pi, P_1) \leftrightarrow \mathcal{L}(A, Ay^1|z)$

$$\bullet dP_0(y|z) \propto e^{-q(y|z)} dP_1(y|z)$$

# Prior Modelling

Full Data  $(A, y', z)$

Observed Data  $(A, A y', z)$

Sensitivity Model  $\text{logit } P(A=1|z, y') = \eta(z) + \underline{q}(y'|z)$

$$\underline{\pi}(z) = P(A=1|z) \quad \underline{p}_i(\cdot|z) = P(y'|A=i, z) \quad P(\cdot|z) = P(y'|z)$$

Strategy 1 prior on  $\beta^z, \underline{\pi}, \underline{p}_i, \underline{q}$ , independent

$$E y' = E_z E(y'|z) = E_z \int y \left[ (1-\underline{\pi}(z)) d\underline{p}_0(y|z) + \underline{\pi}(z) d\underline{p}_1(y|z) \right]$$

$\uparrow \propto e^{\underline{q}(y|z)} d\underline{p}_1(y|z)$

Strategy 2 prior on  $\beta^z, \eta, P, \underline{q}$ , independent

$$E y' = E_z E(y'|z) = E_z \int y dP(y|z)$$

# Findings

Both strategies : posterior of  $Ey'$   $\rightarrow \hat{\theta}$  ,  
even if  $|DATA| \rightarrow \infty$

Strategy 1 : posterior of  $q$  = prior of  $q$

Strategy 2 : posterior of  $q$  depends on DATA

# Nonparametric Bayes (2): Dirichlet process

$$P \sim \text{DP}(a) \iff (P(A_1), \dots, P(A_k)) \sim \text{Dirichlet}(a(A_1), \dots, a(A_k)) \quad \forall \mathcal{X} = \cup_i A_i$$
$$\iff P \sim \sum_{i=1}^{\infty} W_i \delta_{\theta_i}, \quad \theta_i \stackrel{\text{iid}}{\sim} \frac{a}{|a|} \perp W_i = V_i \prod_{j \neq i} (1 - V_j), \quad V_i \stackrel{\text{iid}}{\sim} \text{Be}(1, |a|)$$

**PROP** If  $P \sim \text{DP}(a)$ ,  $X_1, \dots, X_n | P \stackrel{\text{iid}}{\sim} P$ , then  $P | X_1, \dots, X_n \sim \text{DP}(a + nP_n)$

Ferguson, 1973  
Lo, 1983

$$\rightarrow \bullet E(P | X_1, \dots, X_n) = P_n \frac{n}{n+|a|} + a \frac{|a|}{n+|a|} \approx P_n.$$

$$\bullet \sqrt{n} (P - P_n) | X_1, \dots, X_n \rightsquigarrow \mathcal{G}_{P_0}$$

↑ empirical measure

Using  $\text{DP}(a)$  prior is Bayesian equivalent of estimating a distribution by  $P_n$  (= MLE)

# Special Case: No Covariates, given $q$

Full Data  $(A, y')$

Observed Data  $(A, Ay')$

$$\pi = P(A=1), \quad P_1 = \mathcal{L}(y' | A=1), \quad P = \mathcal{L}(y), \quad \text{logit } P(A=1|y) = \eta + q(y)$$

Strategy 1  $\pi \sim \text{Beta}(\alpha, \beta) \perp\!\!\!\perp P_1 \sim \text{DP}(a) \perp\!\!\!\perp q$

Strategy 2  $\eta \sim \mathcal{U}[\alpha, \beta] \perp\!\!\!\perp P \sim \text{DP}(a) \perp\!\!\!\perp q$

**THM** For both strategies, Bernstein-von Mises:

$$\sqrt{n} \left( \underbrace{\theta}_{E y'} - \hat{\theta}_{n, q} \right) \mid \text{DATA}_{n, q} \rightsquigarrow N(0, \sigma_q^2)$$

efficient ("double robust") estimator of  $E y'$



# Proofs

Strategy 1  $\pi \sim \text{Beta}(\alpha, \beta) \perp\!\!\!\perp P_i \sim \text{DP}(a) \perp\!\!\!\perp \mathbf{y}$

$P_i$  is law of data  $\rightarrow$  conjugacy of DP

Strategy 2  $\eta \sim \mathcal{U}[\alpha, \beta] \perp\!\!\!\perp P \sim \text{DP}(a) \perp\!\!\!\perp \mathbf{y}$

$$P \sim \text{DP}(a) \Rightarrow dP_i(y) \propto \frac{1}{1 + e^{\eta + g(y)}} dP(y)$$

$$\Rightarrow P_i \sim \text{"NCRM"}$$

# Nonparametric Bayes (3): NCRM

Kingman, 1975

$\Psi$  Completely Random Measure if  $\Psi(A_1) \dots \Psi(A_k) \Psi(X) \forall X = \bigsqcup_{i=1}^k A_i$

• intensity measure  $\nu = \nu^d + \nu^c$

•  $\Psi = \sum_{i=1}^{\infty} V_i \delta_{a_i} + \sum_{i=1}^{\infty} W_i \delta_{\theta_i}$ ,  $(a_i) \subset X$ ,  $(V_i) \sim \nu^d$ ,  $(W_i, \theta_i) \sim \text{Poisson}(\nu^c)$

$P = \frac{\Psi}{\Psi(X)}$  Normalised C R M

**THM** If  $\nu^c(dy, ds) = s^{-1} e^{-sb(y)} ds d\alpha(y)$ , then for Donsker class  $\mathcal{F}$

$$\sqrt{n} (P - \mathbb{P}_n) | Y_1, \dots, Y_n \rightsquigarrow \mathbb{B}_{P_0} \quad \text{in } \mathcal{L}^p(\mathcal{F})$$

↑ empirical measure

↑ Brownian bridge

$b$  bounded below

$$P_0 b^2 + P_0 \mathcal{F}^r + \alpha \mathcal{F} < \infty$$

$$\frac{1}{q} + \frac{1}{r} < \frac{1}{2}$$

# Nonparametric Bayes (3): NCRM

Kingman, 1975

$\Psi$  Completely Random Measure if  $\Psi(A_1) \perp \dots \perp \Psi(A_k) \quad \forall \mathcal{X} = \bigsqcup_{i=1}^k A_i$

• intensity measure  $\nu = \nu^d + \nu^c$

•  $\Psi = \sum_{i=1}^{\infty} V_i \delta_{a_i} + \sum_{i=1}^{\infty} W_i \delta_{b_i}, \quad (a_i) \subset \mathcal{X}, (V_i) \sim \nu^d, (W_i, b_i) \sim \text{Poisson}(\nu^c)$

$\mathbb{P} = \frac{\Psi}{\Psi(\mathcal{X})}$  Normalised C R M

**THM** If  $\nu^c(dy, ds) = s^{-1-\sigma} e^{-sb(y)} ds d\alpha(y)$ , then for Donsker class  $\mathcal{F}$

$\sqrt{n} (\mathbb{P} - \mathbb{P}_n - \frac{\sigma K_n(G - \hat{\mathbb{P}}_n)}{n}) \mid y_1, \dots, y_n \rightsquigarrow \text{IB}_{\mathbb{P}_0}$  in  $\ell^0(\mathcal{F})$

$\uparrow$  empirical measure  
 $\uparrow$  # distinct values  
 $\uparrow$  empirical measure distinct values  
 $\uparrow$  Gaussian process

$b$  bounded below  
 $\mathbb{P}_0 b^q + \mathbb{P}_0 \mathcal{F}^r + \alpha \mathcal{F} < \infty$

$$\frac{1}{q} + \frac{1}{r} < \frac{1}{2}$$

# proof

**PROP** If  $P \sim \text{NCRM}$ ,  $y^1 \dots y^n | P \stackrel{\text{iid}}{\sim} P$ , then  $P | y_1 \dots y_n \sim \text{mixture of NCRMs}$

James, Lijoi, Prunster, 2009

$P | y_1 \dots y_n, \lambda$  has  $v = v_{n,\lambda}^d + v_{n,\lambda}^c$

$$v_{n,\lambda}^d(\tilde{y}_d, s) \propto s^{N_{y_n}-1} e^{-\lambda(\lambda + b/\tilde{y}_d)} ds$$

↑ distinct values      multiplicity  $\tilde{y}_d$

• Show  $\Pi_n(\lambda \gtrsim \frac{n}{\log n} | y_1 \dots y_n) \xrightarrow{P} 1$

• Show discrete part dominates

(if  $\sigma > 0$  latter is not true)

# Special Case: No Covariates

,  $q \sim \text{prior}$

Full Data  $(A, y')$

Observed Data  $(A, Ay')$

$$\pi = P(A=1), \quad P_1 = \mathcal{L}(y' | A=1), \quad P = \mathcal{L}(y), \quad \text{logit } P(A=1|y) = \eta + \underline{q}(y)$$

Strategy 1  $\quad \pi \sim \text{Beta}(\alpha, \beta) \perp\!\!\!\perp P_1 \sim \text{DP}(\alpha) \perp\!\!\!\perp \underline{q}$

Strategy 2  $\quad \eta \sim \mathcal{U}[\alpha, \beta] \perp\!\!\!\perp P \sim \text{DP}(\alpha) \perp\!\!\!\perp \underline{q}$

**THM** Strategy 1 :  $\underline{q} | \text{DATA} \sim \underline{q} \quad \mathcal{L}(A, Ay')$

Strategy 2 :  $\underline{q} | \text{DATA} \sim \underline{q} | H = H_a$

This works through in posterior of  $E y'$ .

# Special Case: No Covariates

computation

Full Data  $(A, y')$

Observed Data  $(A, Ay')$

$$\pi = P(A=1), \quad P_i = \mathcal{L}(y'_i | A=1), \quad P = \mathcal{L}(y), \quad \text{logit } P(A=1|y) = \eta + \mathbf{q}(y)$$

Strategy 1  $\pi \sim \text{Beta}(\alpha, \beta) \perp\!\!\!\perp P_i \sim \text{DP}(a) \perp\!\!\!\perp \mathbf{q}$

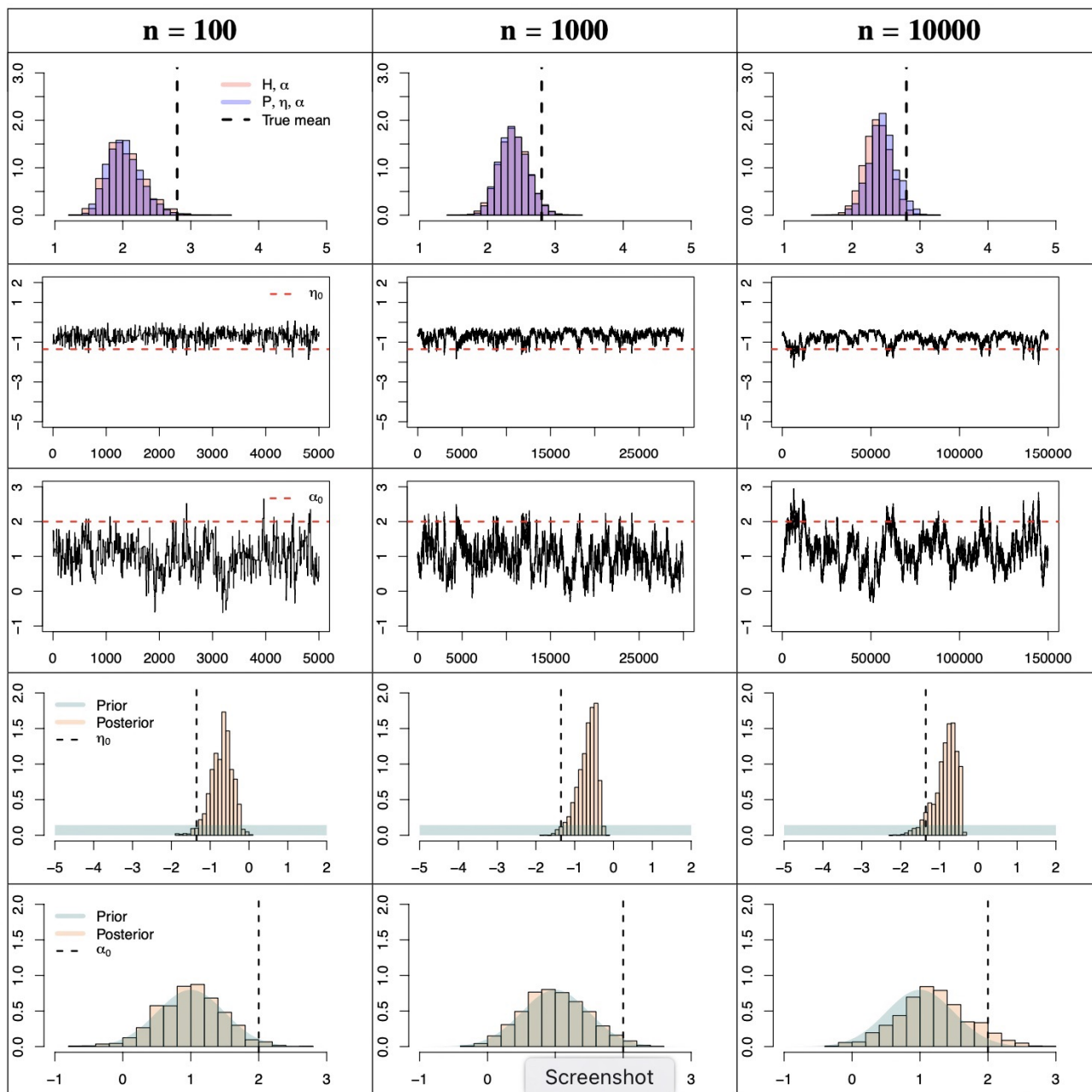
Explicit Dirichlet posterior

Strategy 2  $\eta \sim \mathcal{U}[\alpha, \beta] \perp\!\!\!\perp P \sim \text{DP}(a) \perp\!\!\!\perp \mathbf{q}$

Posterior by Gibbs sampling missing outcomes

# Setting where sensitivity model is not centered at truth

Sensitivity function  $g(y) = \alpha \log y$ ,  $y \in (0, \infty)$



$Eg'$

$\eta$

$\alpha$

$\eta$

$\alpha$

# Special Case: Survival Outcome

Full Data  $(A, y', z)$

$$y' \in (0, \infty)$$

$$P(\cdot | z) = L(y' | z)$$

Observed Data  $(A, Ay', z)$

$$\text{logit } P(A=1 | z, y') = \eta(z) + \underline{g}(y' | z)$$

## Strategy 2

$$P^z \sim \text{DP}(a) \perp\!\!\!\perp \eta(z) = \delta^T z, \delta \sim \Pi \perp\!\!\!\perp P(\cdot | z) \sim \text{Bayesian Cox}$$



# Nonparametric Bayes (4): Cox Model

Cox model  $\Lambda(t|z) = e^{\beta^T z} \Lambda(t)$   
 $\Leftrightarrow P(T > t | z, \beta, \Lambda) = e^{-e^{\beta^T z} \Lambda(t)}$

Prior model  $\beta \sim \pi$ ,  $\Lambda \sim$  Beta process

Hjort, 1990

Data model  $T_1, \dots, T_n | \beta, \Lambda, Z_1, \dots, Z_n \stackrel{\text{ind}}{\sim} \text{Cox}(\beta, \Lambda)$

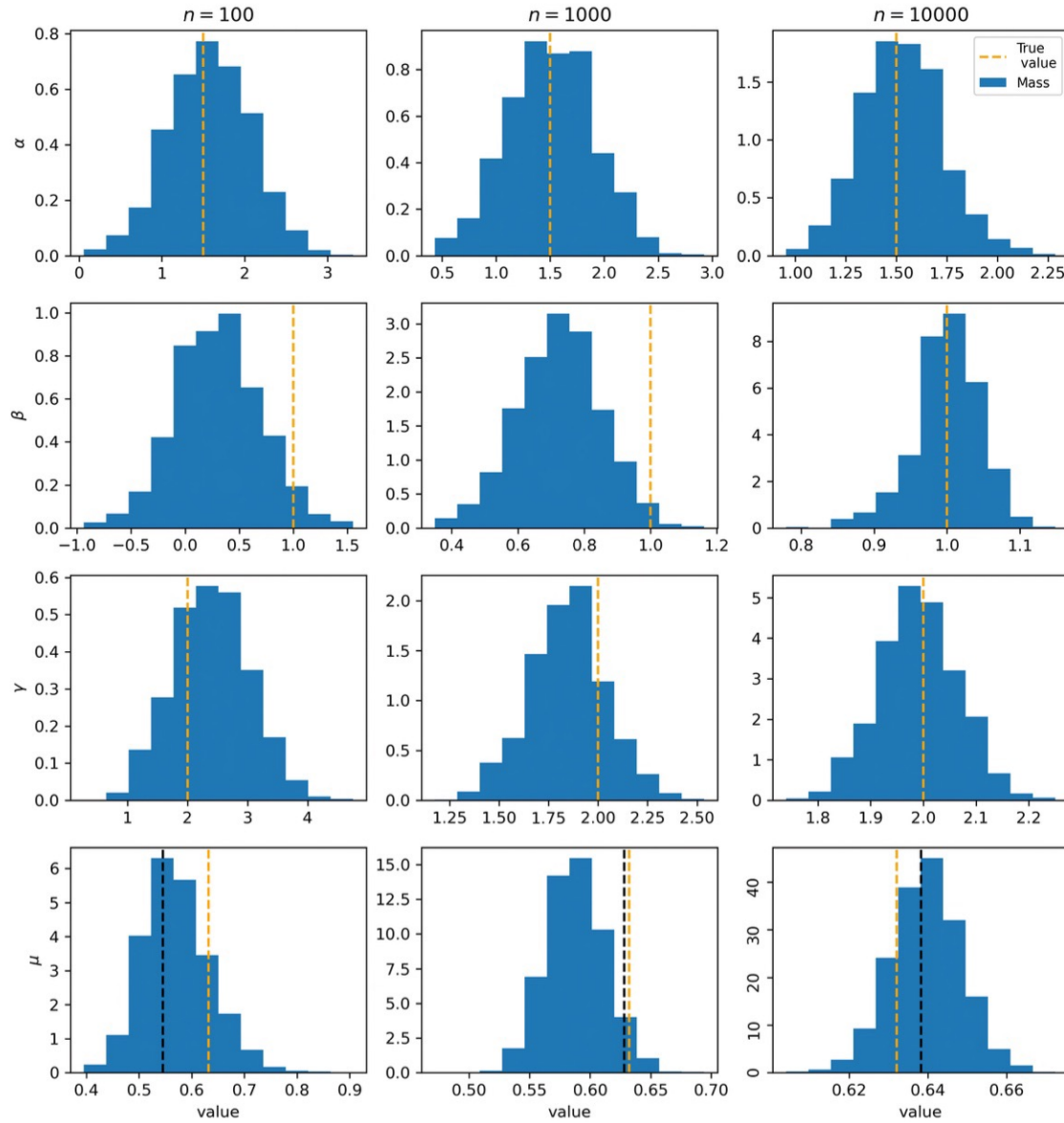
Using Beta process prior is Bayesian equivalent of using the Cox estimator (=MLE)

THM  $\sqrt{n}(\beta - \hat{\beta}_{\text{Cox}}) | T_1, \dots, T_n, Z_1, \dots, Z_n \rightsquigarrow N(0, \tilde{I}_{\text{Cox}}^{-1})$   
 $\sqrt{n}(\Lambda - \hat{\Lambda}_{\text{Cox}}) | T_1, \dots, T_n, Z_1, \dots, Z_n \rightsquigarrow B_{\text{Cox}}$

Kim, Lee, 2003

$$\log_{10} P(A=1 | y', z) = \gamma z + \alpha y' \underbrace{q(y'|z)}_{q_{(0.5, 1)}(z)}$$

$$\Lambda(y'|z) = e^{\beta z} \Lambda(y')$$



$\alpha$

$\beta$

$\gamma$

$E y'$

# Special Case: Binary Outcome

Full Data  $(A, y', z)$

$y' \in \{0, 1\}$

$P(z) = P(y'=1|z)$ ,

Observed Data  $(A, Ay', z)$

$\text{logit } P(A=1|z, y') = \eta(z) + \underline{q}(y'|z)$

Strategy 2  $P^z \sim DP(a) \perp\!\!\!\perp \text{logit } P(y'=1|z) \sim GP \perp\!\!\!\perp \eta \sim GP$   
or Double Robust

Gaussian Process

## CONJECTURE

- $\sqrt{n}(\theta - \hat{\theta}_{n,q}) | \text{DATA}_n, q \rightsquigarrow N(0, \sigma_q^2)$
- $q | \text{DATA}_n$  depends on  $\text{DATA}_n$

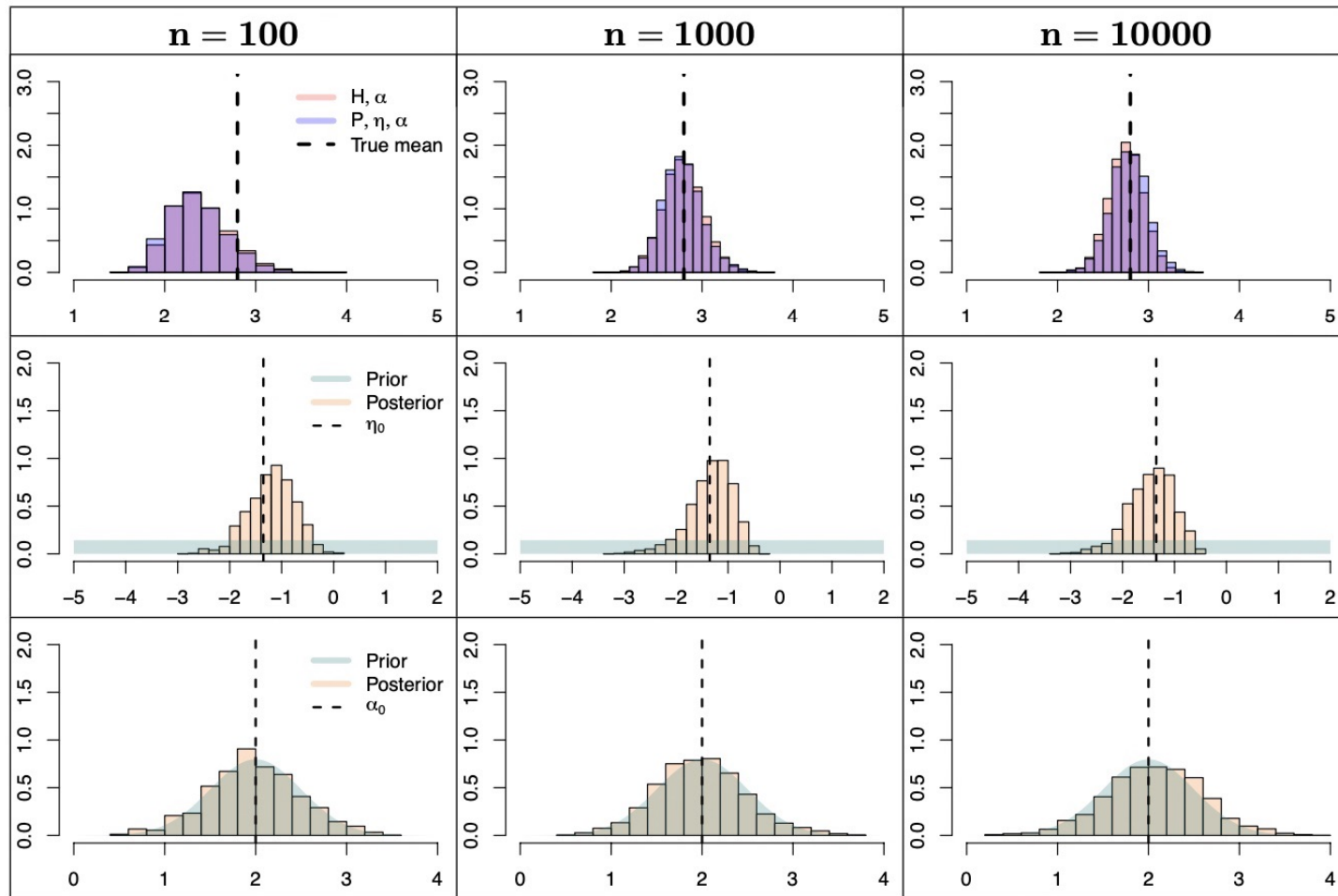
# Conclusion

Even though sensitivity parameter is not identified, in a Bayesian setting part of it can be learned from data, even with nonparametric modelling, depending on prior parameterization.

extra slides →

# Setting where sensitivity model is centered at truth

Sensitivity function  $g(y) = \alpha \log y$ ,  $y \in (0, \infty)$



$E y'$

$\eta$

$\alpha$

# Gibbs Sampling

With Strategy 2 the DP(a) prior is placed on  $P = \mathcal{L}(y_i')$ .  
Posterior given  $(y_i' : i=1, \dots, n)$  is  $DP(a + \sum_{i=1}^n \delta_{y_i'})$ .

We observe  $(y_i' : A_i=1)$  but not  $(y_i' : A_i=0)$ .

→ Generate  $(y_i' : A_i=0)$  using

$$dP_0(y) \propto e^{-q(y)} dP_1(y)$$

and current posterior draws of parameters  $(q, P_1)$ .

→ Repeat

# Future Research : General Outcome

Challenge : nonparametric prior on  $\alpha(y|z)$

- e.g.
- Dependent Dirichlet
  - BART + Gaussian error

Challenge : double robustness (or "debiased machine learning")









