

# Bayesian nonparametric models for treatment effect heterogeneity: model parameterization, prior choice, and posterior summarization

---

Jared S. Murray – The University of Texas at Austin.

March 10, 2022

Bayesian nonparametric modeling is an effective tool for inferring heterogeneous causal effects.

Bayes estimates from these models can have excellent frequentist properties – no need to drink the Kool-Aid.

Some insights about model and prior specification apply to flexible estimation of effect heterogeneity more generally

# Putting BNP to work for inference about effect heterogeneity

Three considerations:

- **Model parameterization:** When you can, isolate your estimand as a parameter
- **Prior specification:** Priors are important for encoding beliefs but also for applying regularization. Regularization that ignores selection can be disastrous.
- **Posterior summarization:** “Solving” the Bayesian analogue of the post-selection inference problem, focusing on stable estimands, and giving actionable insights from complex models.

## Some generic identifying assumptions

*Strong ignorability:*

$$Y_i(0), Y_i(1) \perp\!\!\!\perp Z_i \mid \mathbf{X}_i = \mathbf{x}_i,$$

*Positivity:*

$$0 < \Pr(Z_i = 1 \mid \mathbf{X}_i = \mathbf{x}_i) < 1$$

for all  $i$ . Then

$$P(Y(z) \mid \mathbf{x}) = P(Y \mid Z = z, \mathbf{x})$$

,

and the conditional average treatment effect (CATE) is

$$\begin{aligned} \tau(\mathbf{x}_i) &:= \mathbb{E}(Y_i(1) - Y_i(0) \mid \mathbf{x}_i) \\ &= \mathbb{E}(Y_i \mid \mathbf{x}_i, Z_i = 1) - \mathbb{E}(Y_i \mid \mathbf{x}_i, Z_i = 0). \end{aligned}$$

# Model Parameterization

---

# Parameterizing Nonparametric Models of Causal Effects

Forget confounding and covariates and consider estimating average treatment effect for a binary treatment in a randomized trial.

A simple model:

$$(Y_i | Z_i = 0) \stackrel{iid}{\sim} N(\mu_0, \sigma^2)$$
$$(Y_i | Z_i = 1) \stackrel{iid}{\sim} N(\mu_1, \sigma^2)$$

where the estimand of interest is  $\tau \equiv \mu_1 - \mu_0$ .

If  $\mu_0, \mu_1 \sim N(\phi_j, \delta_j)$  independently then  $\tau \sim N(\phi_1 - \phi_0, \delta_0 + \delta_1)$

Often we have stronger prior information about  $\tau$  than  $\mu_1$  or  $\mu_0$  – in particular, we expect it to be small.

# Parameterizing Nonparametric Models of Causal Effects

A more natural parameterization:

$$(Y_i | Z_i = 0) \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

$$(Y_i | Z_i = 1) \stackrel{iid}{\sim} N(\mu + \tau, \sigma^2)$$

where the estimand of interest is still  $\tau$ .

Now we can express prior beliefs on  $\tau$  directly and independent of nuisance parameters.

# Parameterizing Nonparametric Models of Causal Effects

How does this relate to models for heterogeneous treatment effects?  
Consider (mostly) separate models for treatment arms:

$$y_i = f_{Z_i}(\mathbf{x}_i) + \epsilon_i \quad \epsilon_i \sim N(0, \sigma^2)$$

$$(Y_i | Z_i = 0, \mathbf{x}_i) \stackrel{iid}{\sim} N(f_0(\mathbf{x}), \sigma^2)$$

$$(Y_i | Z_i = 1, \mathbf{x}_i) \stackrel{iid}{\sim} N(f_1(\mathbf{x}), \sigma^2)$$

Independent priors on  $f_0, f_1 \rightarrow$  prior on  $\tau(\mathbf{x}) \equiv f_1(\mathbf{x}) - f_0(\mathbf{x})$  has larger variance than prior on  $f_0$  or  $f_1$

No direct prior control  $\rightarrow$  simple  $f_0, f_1$  can compose to complex  $\tau$  (e.g. Künzel et al (2019)).

In addition, every variable in  $\mathbf{x}$  is a potential effect modifier.



What about the “just another covariate” parameterization?

$$y_i = f(\mathbf{x}_i, z_i) + \epsilon_i \quad \epsilon_i \sim N(0, \sigma^2)$$
$$(Y_i \mid Z_i = z_i, \mathbf{x}_i) \stackrel{iid}{\sim} N(f(\mathbf{x}_i, z_i), \sigma^2)$$

Then the heterogeneous treatment effects are given by

$$\tau(\mathbf{x}) \equiv f(\mathbf{x}, 1) - f(\mathbf{x}, 0)$$

and we still (generally) have no direct prior control!

# Parameterizing Nonparametric Models of Causal Effects

For binary treatments, set  $f(\mathbf{x}_i, z_i) = \mu(\mathbf{x}_i) + \tau(\mathbf{w}_i)z_i$ , where  $\mathbf{w}$  is (possibly) a subset of  $\mathbf{x}$ :

$$y_i = \mu(\mathbf{x}_i) + \tau(\mathbf{w}_i)z_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$

$$(Y_i | Z_i = z_i) \stackrel{iid}{\sim} N(\mu(\mathbf{x}_i) + \tau(\mathbf{w}_i)z_i, \sigma^2)$$

The heterogeneous treatment effects are given by  $\tau(\mathbf{w})$  so we have direct prior control!

In Hahn et. al. (2020), we use independent BART priors on  $\mu$  and  $\tau$  (“Bayesian causal forests”).

## Prior Selection

---

## Tweaking priors on $\tau$

Several adjustments to the BART prior on  $\tau$  in BCF:

- Higher probability on smaller  $\tau$  trees (than BART defaults)
- Higher probability on “stumps” (all stumps = homogeneous effects)
- $N^+(0, \nu)$  Hyperprior on the scale of leaf parameters in  $\tau$

Other nonparametric priors for  $\tau$  have similar “knobs” (scale, smoothness, sparsity, etc.)

For observational data, we need to adjust the prior on  $\mu(\mathbf{x})$  as well, to avoid *regularization induced confounding* (Hahn et al (2016, 2020))

## Regularization can induce confounding (bias)

Let's return to a linear model with homogeneous effects:

$$\begin{aligned}y_i &= f(\mathbf{x}_i, z_i) + \varepsilon_i \\ &= \tau z_i + \beta^t \mathbf{x}_i + \varepsilon_i\end{aligned}$$

and suppose  $x_i$  is high dimensional.

Assume  $\beta \sim N(0, \lambda^{-1}I)$  (ridge prior) and  $p(\tau) \propto 1$

What effect does the prior (regularization) have on estimating  $\tau$  using the posterior mean?

## Regularization can induce confounding (bias)

The bias of  $\tilde{\tau} = E(\tau | Y, z, \mathbf{x})$  is

$$\text{bias}(\tilde{\tau}) = \lambda \hat{\delta}^t [\lambda \mathbf{I} + \mathbf{X}^t (\mathbf{I} - \mathbf{P}_z) \mathbf{X}]^{-1} \beta \quad (1)$$

where  $\hat{\delta}_j$  = the OLS estimate of  $x_{ij} = \delta_j z_i + \epsilon_{ij}$ . Alternatively:

$$\text{bias}(\tilde{\tau}) = \lambda [\mathbf{z}^t (\mathbf{z} - \tilde{\mathbf{z}}_\lambda)]^{-1} \tilde{\gamma}_\lambda^t \beta \quad (2)$$

where  $\tilde{\gamma}_\lambda = [\lambda \mathbf{I} + \mathbf{X}^t \mathbf{X}]^{-1} \mathbf{X}^t \mathbf{z}$  and  $\tilde{\mathbf{z}}_\lambda = \mathbf{X} \tilde{\gamma}_\lambda$

In general, if  $z$  and  $\mathbf{x}$  are correlated the bias is nonzero and depends on the nuisance parameter!

## Solution: Don't penalize variation in $f(\mathbf{x}, z)$ along $E(Z | \mathbf{x})$

Expand the model to include  $\hat{z}_i$  (a function of  $z$  and  $\mathbf{X}$ ) that estimates  $E(Z | \mathbf{x})$ :

$$\begin{aligned}y_i &= f(\mathbf{x}_i, z_i) + \varepsilon_i \\ &= \tau z_i + \phi \hat{z}_i + \beta^t \mathbf{x}_i + \varepsilon_i\end{aligned}$$

Keep  $\beta \sim N(0, \lambda^{-1}I)$  (ridge prior) with  $p(\tau, \phi) \propto 1$ , so that variation in the direction of  $\hat{z}_i$  is unregularized

$$\text{bias}(\tilde{\tau}) = \lambda \hat{\delta}^t [\lambda I + \mathbf{X}^t (\mathbf{I} - \mathbf{P}_z) \mathbf{X}]^{-1} \beta \quad (3)$$

where  $\hat{\delta}_j =$  the OLS estimate of  $x_{ij} = \alpha_j \hat{z}_i + \delta_j z_i + \epsilon_{ij}$  ( $\approx 0$ ).

## Solution: Don't penalize variation in $f(\mathbf{x}, z)$ along $E(Z | \mathbf{x})$

Expand the model to include  $\hat{z}_i$  (a function of  $z$  and  $\mathbf{X}$ ) that estimates  $E(Z | \mathbf{x})$ :

$$\begin{aligned}y_i &= f(\mathbf{x}_i, z_i) + \varepsilon_i \\ &= \tau z_i + \phi \hat{z}_i + \beta^t \mathbf{x}_i + \varepsilon_i\end{aligned}$$

Keep  $\beta \sim N(0, \lambda^{-1}I)$  (ridge prior) with  $p(\tau, \phi) \propto 1$ , so that variation in the direction of  $\hat{z}_i$  is unregularized

$$\text{bias}(\tilde{\tau}) = \lambda \hat{\delta}^t [\lambda I + \mathbf{X}^t (\mathbf{I} - \mathbf{P}_z) \mathbf{X}]^{-1} \beta \quad (4)$$

where  $\hat{\delta}_j =$  the OLS estimate of  $x_{ij} = \alpha_j \hat{z}_i + \delta_j z_i + \varepsilon_{ij}$  ( $\approx 0$ ).



# Regularization induced confounding is a general phenomenon

There is nothing special about the ridge prior or the linear model – RIC is easy to produce with nonlinear models and nonlinear data generating processes. (Hahn et al (2020))

In essence: Since  $Z$  is a proxy for  $E(Z | \mathbf{x})$ , if the prior on  $f(\mathbf{x}, z)$  strongly penalizes variation in the “direction” of  $E(Z | \mathbf{x})$  (and not  $Z$ ) the prior encourages misattributing that variation in  $f$  to  $Z$ .

This is not a Bayes problem; it’s a generic regularization problem.

## How to avoid penalizing variation in $f(\mathbf{x}, z)$ along $E(Z | \mathbf{x})$

Including  $\hat{z}_i$  as an extra coordinate/feature/covariate is often enough to mitigate regularization induced confounding.

Depending on the model, there may be easier/more efficient ways to accomplish this (e.g. residualization).

In Hahn et al (2020) we evaluate BART priors on  $f(\mathbf{x}, z)$  with and without  $\hat{z}$  and BCF:

$$y_i = \mu(\mathbf{x}_i, \hat{z}_i) + \tau(\mathbf{w}_i)z_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$

The latter two are often *much* better and rarely worse, especially when selection into treatment is based on expected outcomes under control ("targeted selection").

# Posterior Summarization

---

# Posterior summaries, or: I fit this model, now what?

Examine the “best” (in a user-defined sense) simple approximation to a “true”  $g(\mathbf{x})$  (Woody et al (2020))

Given samples of a function  $g(\mathbf{x})$ ,

1. Consider a class of simple/interpretable approximations  $\Gamma$  to  $g$
2. Make inference on

$$\gamma = \arg \min_{\tilde{\gamma} \in \Gamma} d(g, \tilde{\gamma}, \tilde{\mathbf{X}}) + p(\tilde{\gamma})$$

for an appropriate distance function  $d$  and (optional) complexity penalty  $p(\gamma)$

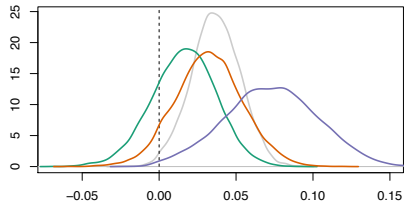
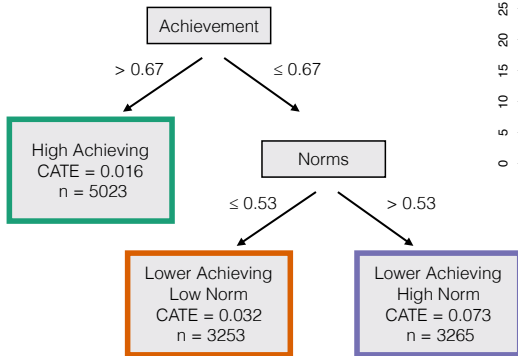
Get *draws* of  $\gamma$  by solving the optimization for each draw of  $g$ . Get point estimates by solving

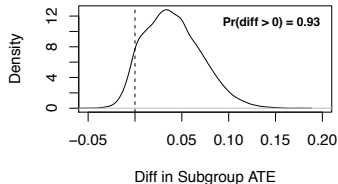
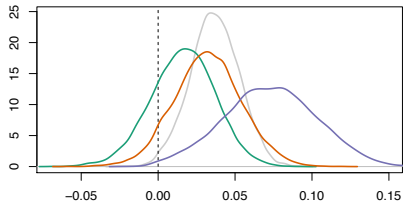
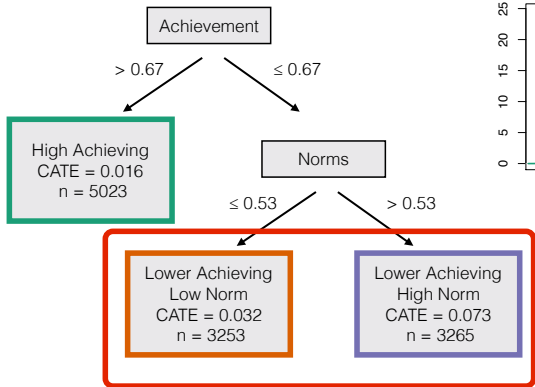
$$\hat{\gamma} = \arg \min_{\tilde{\gamma} \in \Gamma} E_g[d(g, \tilde{\gamma}, \tilde{\mathbf{X}}) + p(\tilde{\gamma}) \mid Y, \mathbf{x}]$$

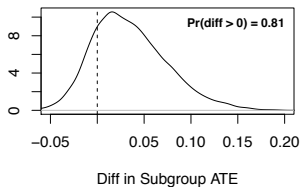
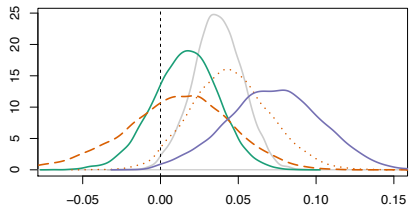
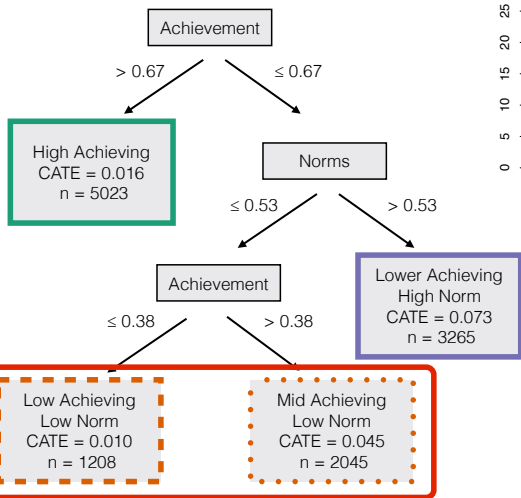
# Posterior summaries, or: I fit this model, now what?

Posterior summaries:

1. Are more interpretable (subgroup analysis, linear/additive/sparse approximations) and can be targeted to scientific questions
2. Obviate the “need” to fit multiple models for different questions (Bayesians need to think about post-selection issues too) – multiple summaries use the data *once* to go prior → posterior
3. Are often more stable (coarse subgroup effects vs. individualized estimates)
4. Come with (Bayes) valid estimates of uncertainty









# Additive summaries

We can get approximate partial effect curves via additive summaries:

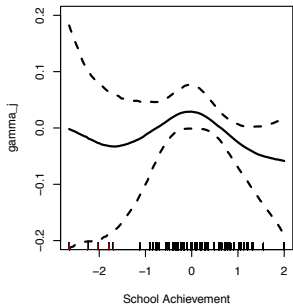
$$\tau(\mathbf{w}) \approx \gamma_0 + \sum_{j=1}^p \gamma_j(w_j)$$

with appropriate forms for  $\gamma_j$  plus smoothing penalties.

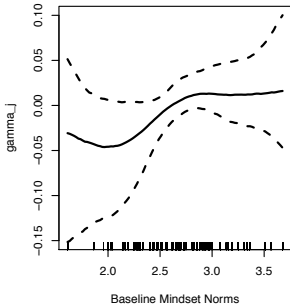
We can also get posterior on discrepancy metrics, like pseudo- $R^2$ :

$$\text{Cor}^2(\gamma(\mathbf{w}_i), \tau(\mathbf{w}_i))$$

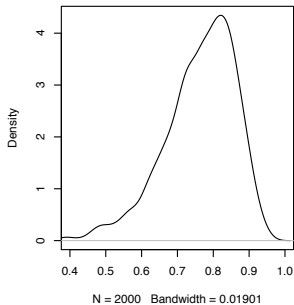
**Approx Additive Effect**



**Approx Additive Effect**

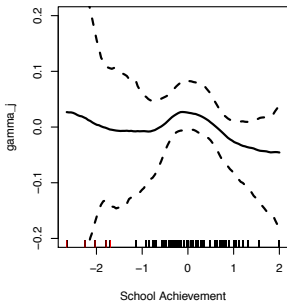


**Approximation R<sup>2</sup>**

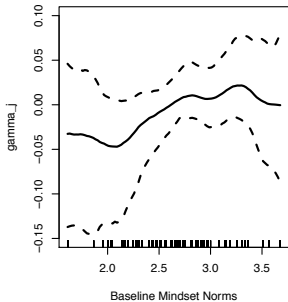


(Partial effect of minority composition not shown)

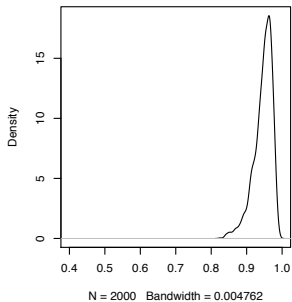
**Approx Additive Effect**



**Approx Additive Effect**



**Approximation R^2**



(Partial effect of minority composition not shown)

## Other applications of posterior summarization

- Interaction detection (Woody et al (2020)): If an additive summary is poor how do we search for missing interactions?
- Sensitivity to control function specification (Woody et al (2020b)): How do I expect removing confounders (or nonlinear/interaction terms) to change my effect estimate?
- “Explanations”: Linear summaries in neighborhoods of  $\mathbf{x}_i$  = LIME with uncertainty

# Thank you!

[jared.murray@mcombs.utexas.edu](mailto:jared.murray@mcombs.utexas.edu)

<https://jaredsmurray.github.io>