

# Recent developments in model specification, regularization, and summarization for nonparametric Bayesian models of heterogeneous treatment effects

---

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

Joint work with Carlos Carvalho, P. Richard Hahn, David Yeager, et al.

June, 2019

# National Study of Learning Mindsets

**MINDSET**  
SCHOLARS  
NETWORK

[About Us](#)[About Mindsets](#)[Research & Resources](#)[Mindsets in the News](#)[Blog](#)

## Mindsets Matter

How might learning environments influence students' mindsets?

### LATEST RESEARCH



**Reducing Racial Gaps in School Suspensions**  
Brief Intervention to Encourage Empathic Discipline Cuts Suspension Rates in Half Among Adolescents  
[READ MORE >](#)



**Mindset Programs That Improve College Outcomes**  
Teaching a Lay Theory Before College Narrows Achievement Gaps at Scale  
[READ MORE >](#)



**Parent Practices & Children's Mindsets**  
What Predicts Children's Fixed and Growth Intelligence Mindsets? Not Their Parents' Views of Intelligence But Their Parents' Views of Failure  
[READ MORE >](#)

# National Study of Learning Mindsets

- National Study of Learning Mindsets (Yeager et. al., 2017):  
Randomized controlled trial of a low-cost mindset intervention
- Probability sample, 65 schools in this analysis (> 11,000 9th grade students)
- Specifically designed to assess treatment effect heterogeneity

# National Study of Learning Mindsets

What do we hope to gain with BNP?

- Avoid (explicit) model selection/specification search
- Flexible models of treatment effect heterogeneity
- Appropriate measures of uncertainty

Where does BNP need a little help?

- Summarizing complex posteriors to communicate results

# Our (generic) 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

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

# Parameterizing Nonparametric Models of Causal Effects

Let's forget confounding and covariates for a second.

A simple model:

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

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

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

If  $\mu_j \sim N(\phi_j, \delta_j)$  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:

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

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

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

Now we can express prior beliefs on  $\tau$  directly.

# Parameterizing Nonparametric Models of Causal Effects

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

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

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

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

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 – simple  $f_0, f_1$  can compose to complex  $\tau$ .

Every variable in  $\mathbf{x}$  is a potential effect modifier.



# Parameterizing Nonparametric Models of Causal Effects

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 given by

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

and every variable in  $\mathbf{x}$  is still a potential effect modifier, and we have no direct prior control...

# Parameterizing Nonparametric Models of Causal Effects

Set  $f(\mathbf{x}_i, z_i) = \mu(\mathbf{x}_i) + \tau(\mathbf{w}_i)z_i$ , where  $\mathbf{w}$  is a subset of  $\mathbf{x}$ :

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

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

The heterogeneous treatment effects are given by  $\tau(\mathbf{w})$  directly

In Hahn et. al. (2017) we use independent BART priors on  $\mu$  and  $\tau$  (“Bayesian causal forests”)

# Tweaking priors for causal effects

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

- Higher probability on smaller  $\tau$  trees (than BART defaults)
- Higher probability on “stumps”, trees that never split (all stumps = homogeneous effects)
- $N^+(0, \nu)$  Hyperprior on the scale of leaf parameters in  $\tau$

# What about observational data?

What changes when adding confounding?

Not much, but we should include an estimated propensity score as a covariate:

$$y_i = \mu(\mathbf{x}_i, \hat{\pi}(\mathbf{x}_i)) + \tau(\mathbf{w}_i)Z_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$
$$(Y_i \mid Z = z_i) \stackrel{iid}{\sim} N(\mu(\mathbf{x}_i, \hat{\pi}(\mathbf{x}_i)) + \tau(\mathbf{w}_i)z_i, \sigma^2)$$

Mitigates regularization induced confounding (RIC); see Hahn et. al. (2017) for details.

Lots of independent empirical evidence this is important (cf Dorie et al (2019), Wendling et al (2019)).

# The National Study of Learning Mindsets

**MINDSET**  
SCHOLARS  
NETWORK

SEARCH

Q


About Us

About Mindsets

Research & Resources

Mindsets in the News

Blog




**Mindsets Matter**


How might learning environments influence students' mindsets?

● ○ ○


LATEST RESEARCH



**Reducing Racial Gaps in School Suspensions**  
Brief Intervention to Encourage Empathic Discipline Cuts Suspension Rates in Half Among Adolescents  
[READ MORE >](#)



**Mindset Programs That Improve College Outcomes**  
Teaching a Lay Theory Before College Narrows Achievement Gaps at Scale  
[READ MORE >](#)



**Parent Practices & Children's Mindsets**  
What Predicts Children's Fixed and Growth Intelligence Mindsets? Not Their Parents' Views of Intelligence But Their Parents' Views of Failure  
[READ MORE >](#)

Answering Questions, Working Toward Solutions

# Analysis of National Study of Learning Mindsets

- A new analysis of effects on math GPA in the overall population of students
- Interesting moderators are baseline level of mindset norms, school achievement, and minority composition
- Many, many other controls

Mindset study has students nested within schools:

$$y_{ij} = \alpha_j + \mu(\mathbf{x}_{ij}) + [\phi_j + \tau(\mathbf{w}_{ij})] z_{ij} + \epsilon_{ij}, \quad \epsilon_{ij} \sim N(0, \sigma^2)$$

$\alpha_j, \phi_j$  have standard random effect priors (normal with half- $t$  hyperpriors on scale)

$\mathbf{w}$  = school achievement, baseline mindset norms, minority composition

# We fit this model, now what?

How do we present our results? Average treatment effect (ATE), school ATE, plots of (school) ATE by  $x_j$ ....

Two questions from our collaborators:

- For which groups was the treatment most/least effective?
- What is the impact of posited treatment effect modifiers, holding other modifiers constant?



# We fit this model, now what?

How do we present our results? Average treatment effect (ATE), school ATE, plots of (school) ATE by  $x_j$ ....

Two questions from our collaborators:

- For which groups was the treatment most/least effective?
- What is the impact of posited treatment effect modifiers, holding other modifiers constant?

# Subgroup finding as a decision problem

The action  $\gamma$  is choosing subgroups, here represented by a recursive partition (tree)

- Minimize the posterior expected loss

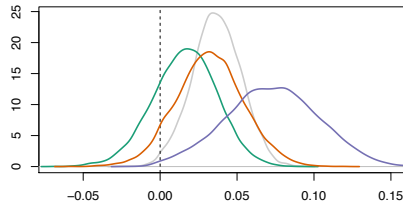
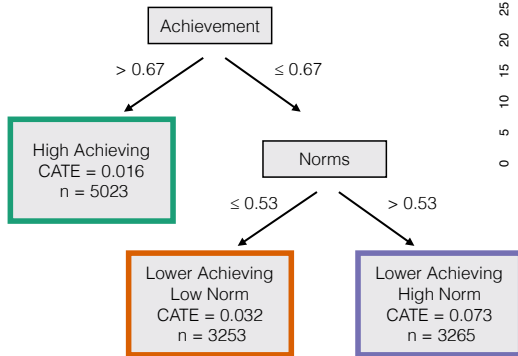
$$\hat{\gamma} = \arg \min_{\tilde{\gamma} \in \Gamma} E_{\tau} [d(\tau, \tilde{\gamma}) + p(\tilde{\gamma}) \mid Y, \mathbf{x}]$$

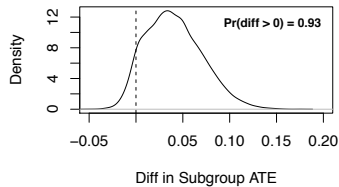
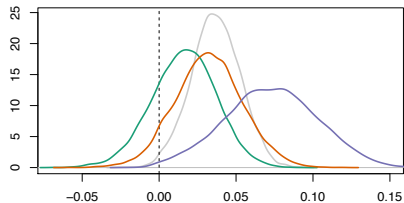
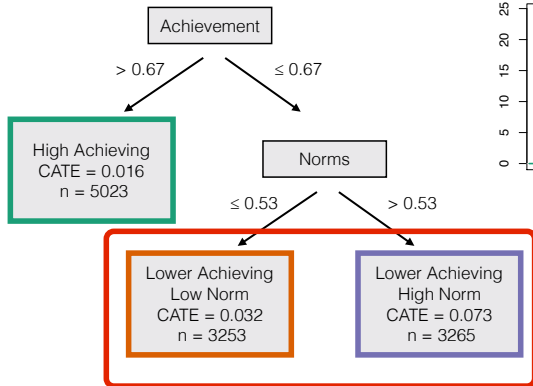
where  $p()$  is a complexity penalty and  $d()$  is squared error averaged over a distribution for  $\mathbf{x}$ .

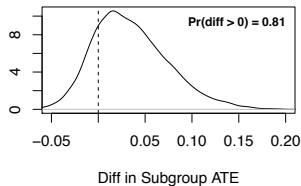
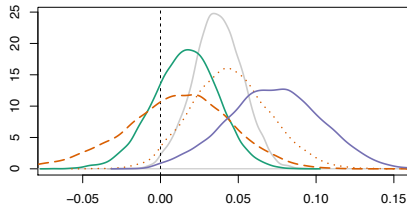
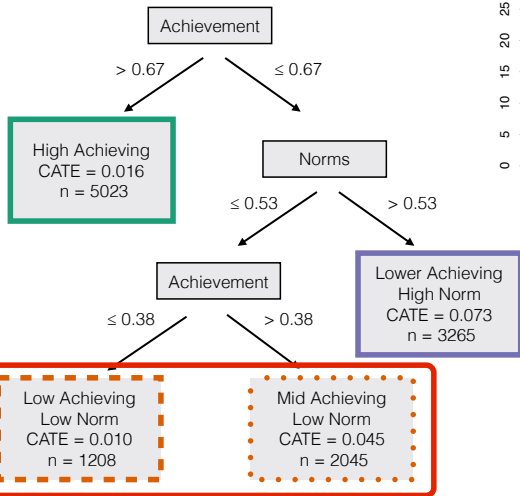
- With  $\hat{\gamma}$  in hand we can look at the **joint** posterior distribution of subgroup ATEs

Relaxing complexity penalties in the loss function  $\rightarrow$  growing a deeper tree

(Hahn et al (2017), Sivaganesan et al (2017))







# We fit this model, now what?

How do we present our results? ATE, school ATE, plots of school ATE by variables....

Two questions we got from our collaborators:

- For which groups was the treatment most/least effective?
- What was the impact of posited treatment effect modifiers, holding other modifiers constant?

We could generate posterior distributions of individual-level partial effects manually and try to summarize these, **or** try to approximate  $\tau$  by a proxy with simple partial effects (i.e. an additive function)

# Posterior projections for model interpretations with uncertainty

The goal: Examine the “best” (in a user-defined sense) simple approximation to the “true”  $\tau(\mathbf{x})$

Given samples of  $\tau$ ,

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

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

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

Get draws of  $\gamma$  by solving the optimization for each draw of  $\tau$

Also get discrepancy metrics, like pseudo- $R^2$ :  $\text{Cor}^2[\gamma(\mathbf{x}_i), \tau(\mathbf{x}_i)]$

(Woody, Carvalho, and Murray (2019) for this approach in predictive models)

# Approximate partial effects of treatment effect modifiers

Here we use:

- An additive approximation  $\gamma(\mathbf{x})$
- $d(\tau, \gamma) = \sum_{i=1}^n (\tau(\mathbf{x}_i) - \gamma(\mathbf{x}_i))^2$
- A smoothness penalty  $p(\gamma)$

Don't average draws to get a point estimate! Treat

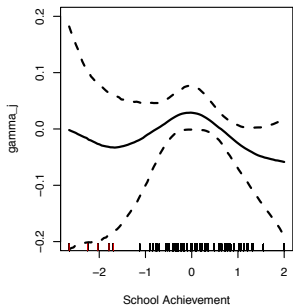
$$d(\tau, \tilde{\gamma}) + p(\tilde{\gamma}) = \sum_{i=1}^n (\tau(\mathbf{x}_i) - \tilde{\gamma}(\mathbf{x}_i))^2 + p(\tilde{\gamma})$$

as a loss function, and minimize posterior expected loss:

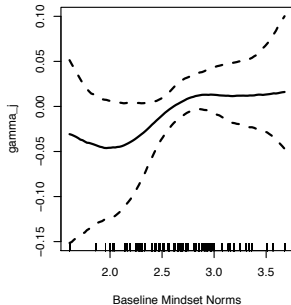
$$\hat{\gamma} = \arg \min_{\tilde{\gamma} \in \Gamma} E_{\tau}[d(\tau, \tilde{\gamma}) + p(\tilde{\gamma}) \mid Y, \mathbf{x}]$$



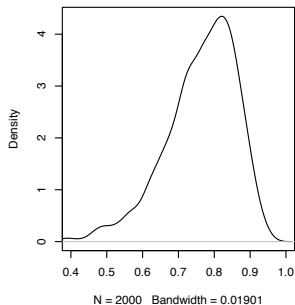
**Approx Additive Effect**



**Approx Additive Effect**

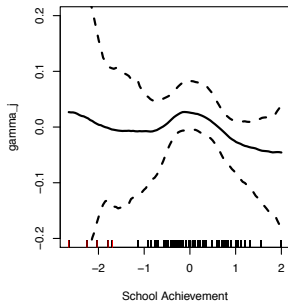


**Approximation  $R^2$**

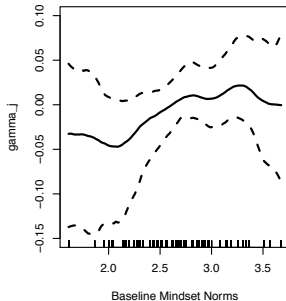


(Partial effect of minority composition not shown)

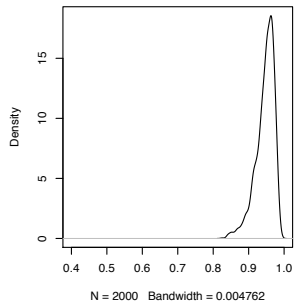
**Approx Additive Effect**



**Approx Additive Effect**



**Approximation R^2**



(Partial effect of minority composition not shown)

# Thank you!

- P. Richard Hahn, Jared S. Murray, Carlos M. Carvalho: “Bayesian regression tree models for causal inference: regularization, confounding, and heterogeneous effects”, 2017; arXiv:1706.09523.
- Spencer Woody, Carlos M. Carvalho, Jared S. Murray: “Model interpretation through lower-dimensional posterior summarization”, 2019; arXiv:1905.07103.