

# Bayesian Inference of the Number of Trees in the BART Model

Gavin Collins<sup>1,2</sup>, Matthew Pratola<sup>2</sup>, Radu Herbei<sup>2</sup>, Robert McCulloch<sup>3</sup>, and Edward George<sup>4</sup>

Brigham Young University

45<sup>th</sup> Annual Summer Institute of Applied Statistics

June 22, 2023

<sup>1</sup>Sandia National Laboratories, <sup>2</sup>The Ohio State University, <sup>3</sup>Arizona State University,

<sup>4</sup>Wharton School at the University of Pennsylvania

# The Bottom Line Up Front

- Prior distribution on the number of trees
- MH step to add/delete one tree at a time
- Takes longer
- Works well
- (Still a work in progress)

1. Recap of BART
2. Bayesian Inference of the Number of Trees
  - i. Motivation
  - ii. A Fully Bayesian Model
  - iii. Sampling from the Posterior Distribution
  - iv. Code
  - v. Simulations
  - vi. Application to Real Data
3. Conclusion

# 1. Recap of BART

## 2. Bayesian Inference of the Number of Trees

i. Motivation

ii. A Fully Bayesian Model

iii. Sampling from the Posterior Distribution

iv. Code

v. Simulations

vi. Application to Real Data

## 3. Conclusion

# Bayesian Additive Regression Trees

## Data

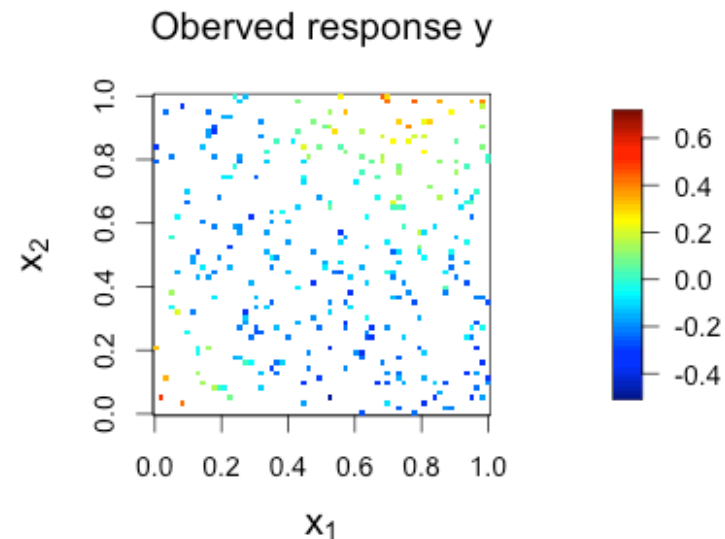
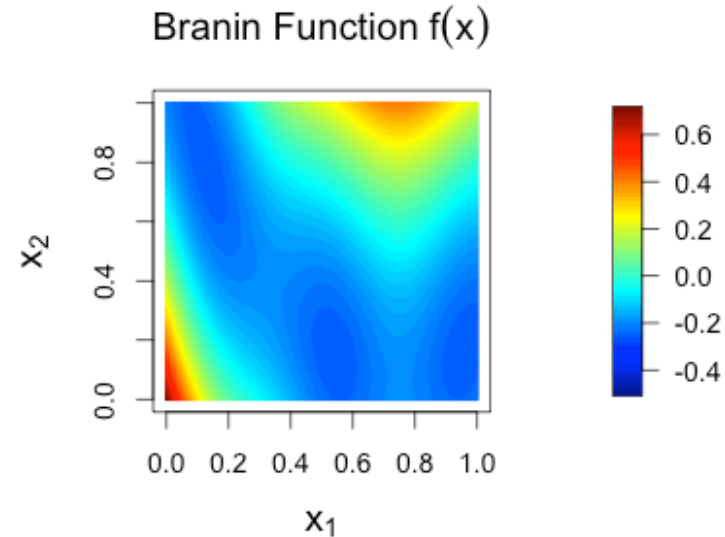
- $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$
- input  $\mathbf{x}_i \in \mathbb{R}^p \rightarrow$  response  $y_i \in \mathbb{R}$

## Regression Model

- $y_i | \mathbf{x}_i \sim N(f(\mathbf{x}_i), \sigma^2), i = 1, \dots, n$  (*iid*)
- $f: \mathbb{R}^p \rightarrow \mathbb{R}$  (mean function)
- $\sigma^2 \geq 0$  (residual variance)

## “Branin” Example:

- $p = 2$
- $f =$  “The Branin Function”
- $n = 300$
- $\mathbf{x}_1, \dots, \mathbf{x}_{300} \sim \text{Unif}(0,1)^2$  (*iid*)
- $\sigma^2 = 1$



# Bayesian Additive Regression Trees

## Data

- $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$
- input  $\mathbf{x}_i \in \mathbb{R}^p \rightarrow$  response  $y_i \in \mathbb{R}$

## Regression Model

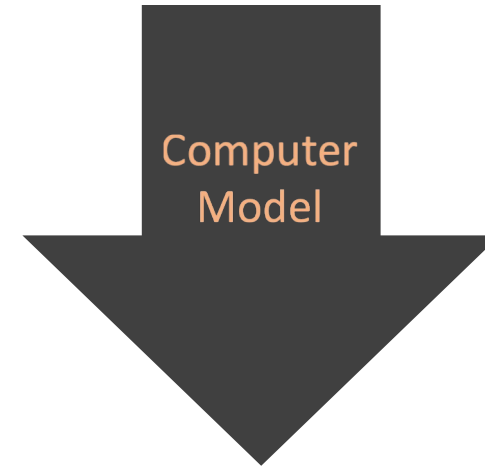
- $y_i | \mathbf{x}_i \sim N(f(\mathbf{x}_i), \sigma^2), i = 1, \dots, n$  (ind)
- $f: \mathbb{R}^p \rightarrow \mathbb{R}$  (mean function)
- $\sigma^2 \geq 0$  (residual variance)

## Hurricane Example:

- $p = 6$
- $f =$  Computer Model
- $n = 4,000$
- Goal: Infer  $f$  for sensitivity analysis, model calibration, etc.

## Input $\mathbf{x}$

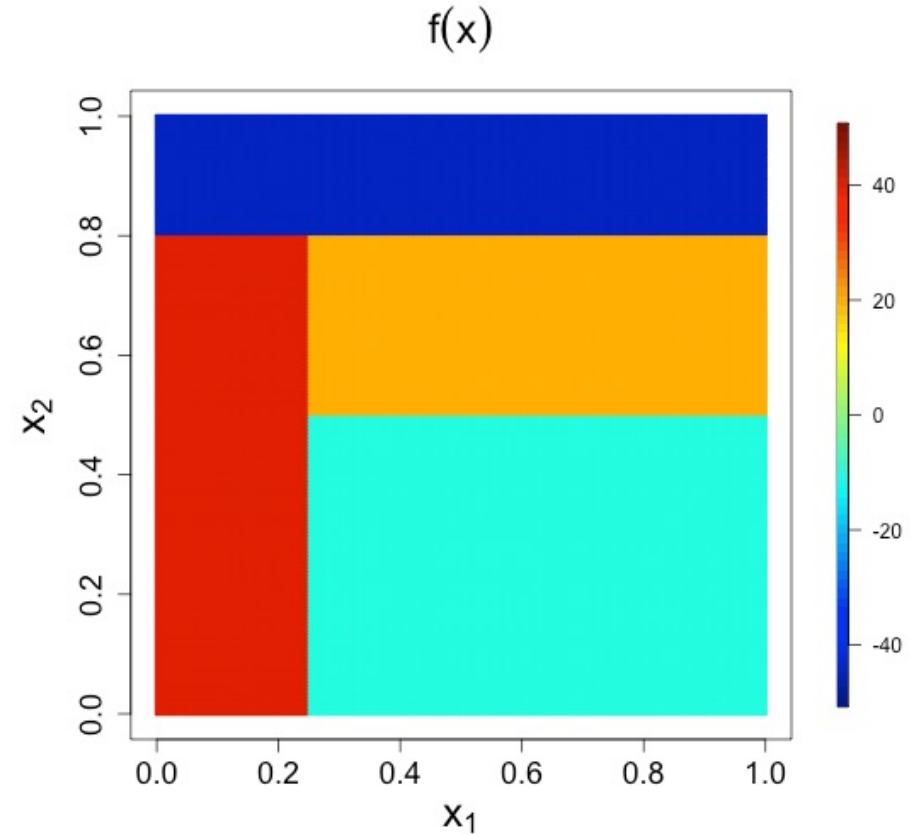
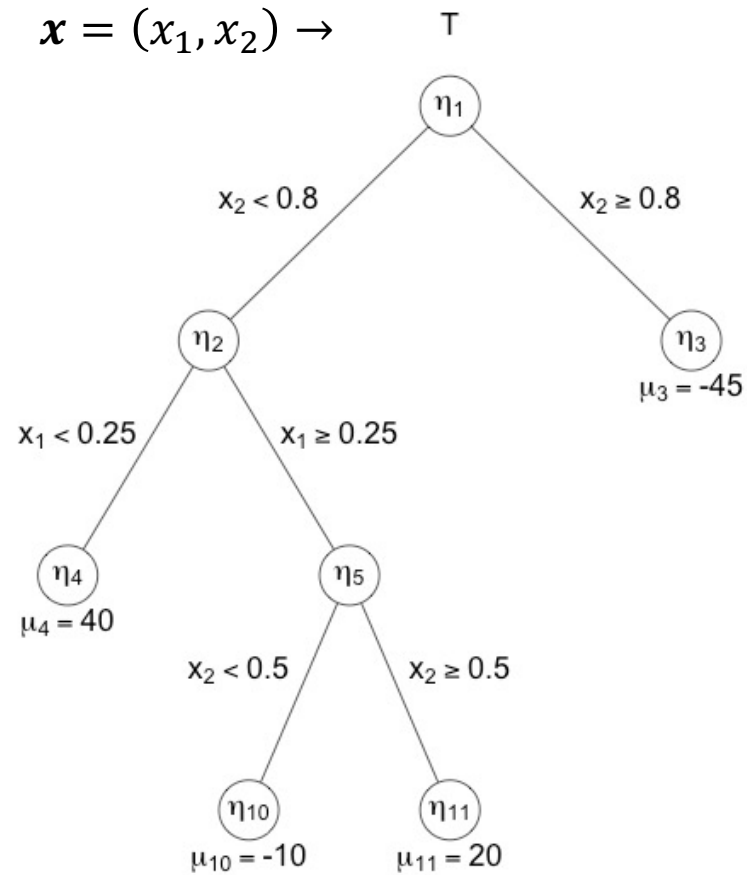
- $x_1 =$  Initial Sea Level
- $x_2 =$  Hurricane Heading
- $x_3 =$  Velocity of the Eye
- $x_4 =$  Max Wind Speed
- $x_5 =$  Min Pressure
- $x_6 =$  Landfall Location



## Response $y$

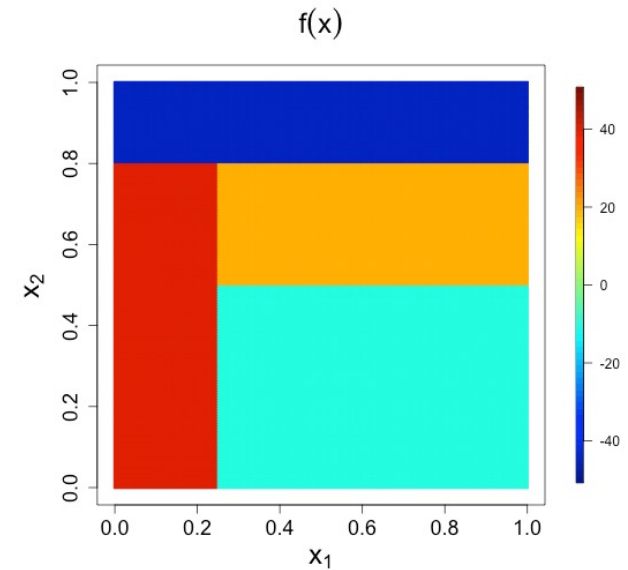
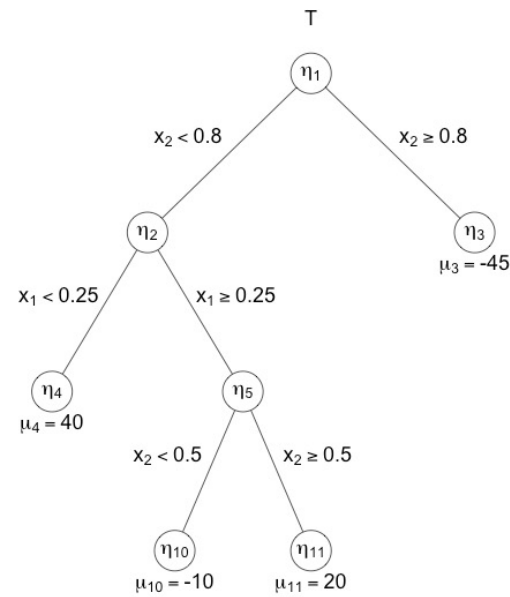
$y =$  Maximum Water Level During a Storm Surge

# Bayesian Additive Regression Trees



# Bayesian Additive Regression Trees

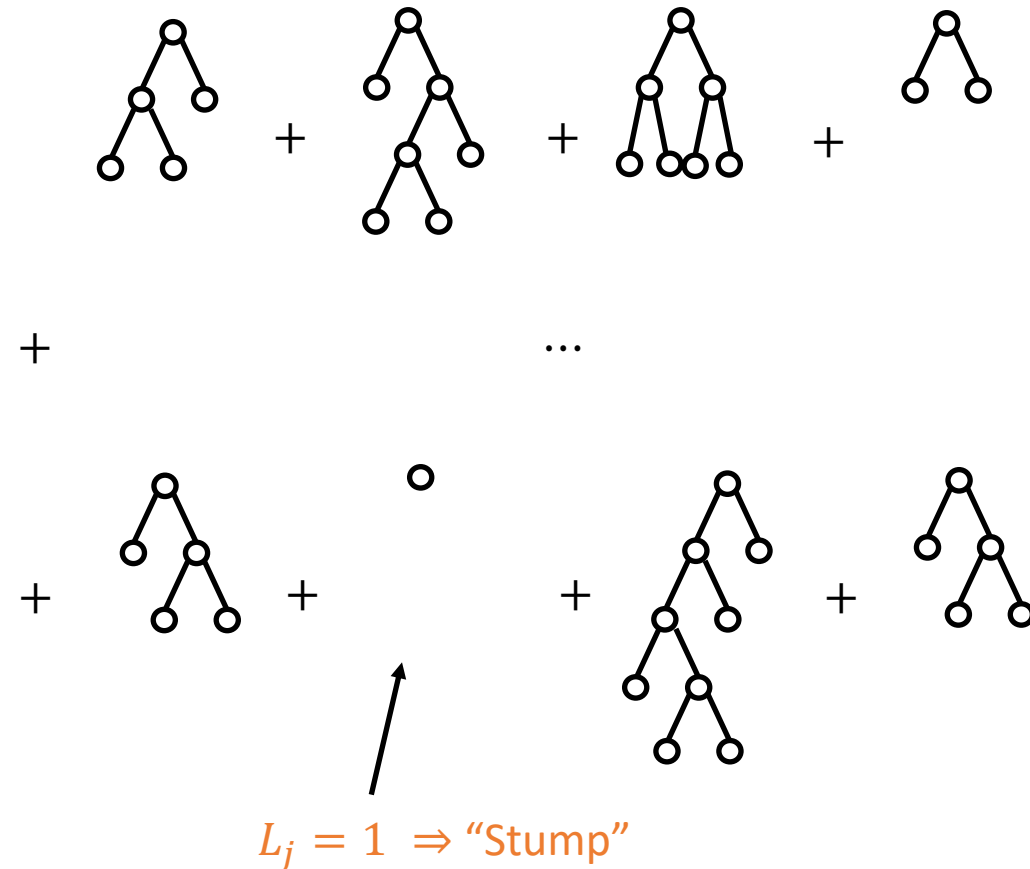
- $T$  has  $L$  “terminal nodes”
- Terminal node parameter  $\mu \in \mathbb{R}^L$
- $f(\mathbf{x}) = g(\mathbf{x}; T, \mu)$





# Bayesian Additive Regression Trees

- $T_{1:m} \equiv T_1, \dots, T_m$  ( $m \approx 200$ )
- $\boldsymbol{\mu}_{1:m} \equiv \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_m$
- $T_j$  has  $L_j$  terminal nodes
- $\boldsymbol{\mu}_j \in \mathbb{R}^{L_j}$
- $f(\mathbf{x}) = \sum_{j=1}^m g(\mathbf{x}; T_j, \boldsymbol{\mu}_j)$



# Bayesian Additive Regression Trees

**Prior Distribution**  $\pi(T_{1:m}, \boldsymbol{\mu}_{1:m}, \sigma^2) = \pi(\sigma^2) \prod_{j=1}^m \pi(T_j) \pi(\boldsymbol{\mu}_j | T_j)$

- $T_j \sim$  Tree-Generating Stochastic Process
- $\mu_{j\ell} | T_j \sim N(0, \tau_m^2)$ ;  $\ell = 1, \dots, L_j$ ;  $j = 1, \dots, m$  (iid)
- $\sigma^2 \sim$  Scaled-inv- $\chi^2(\nu, \lambda)$

# Bayesian Additive Regression Trees

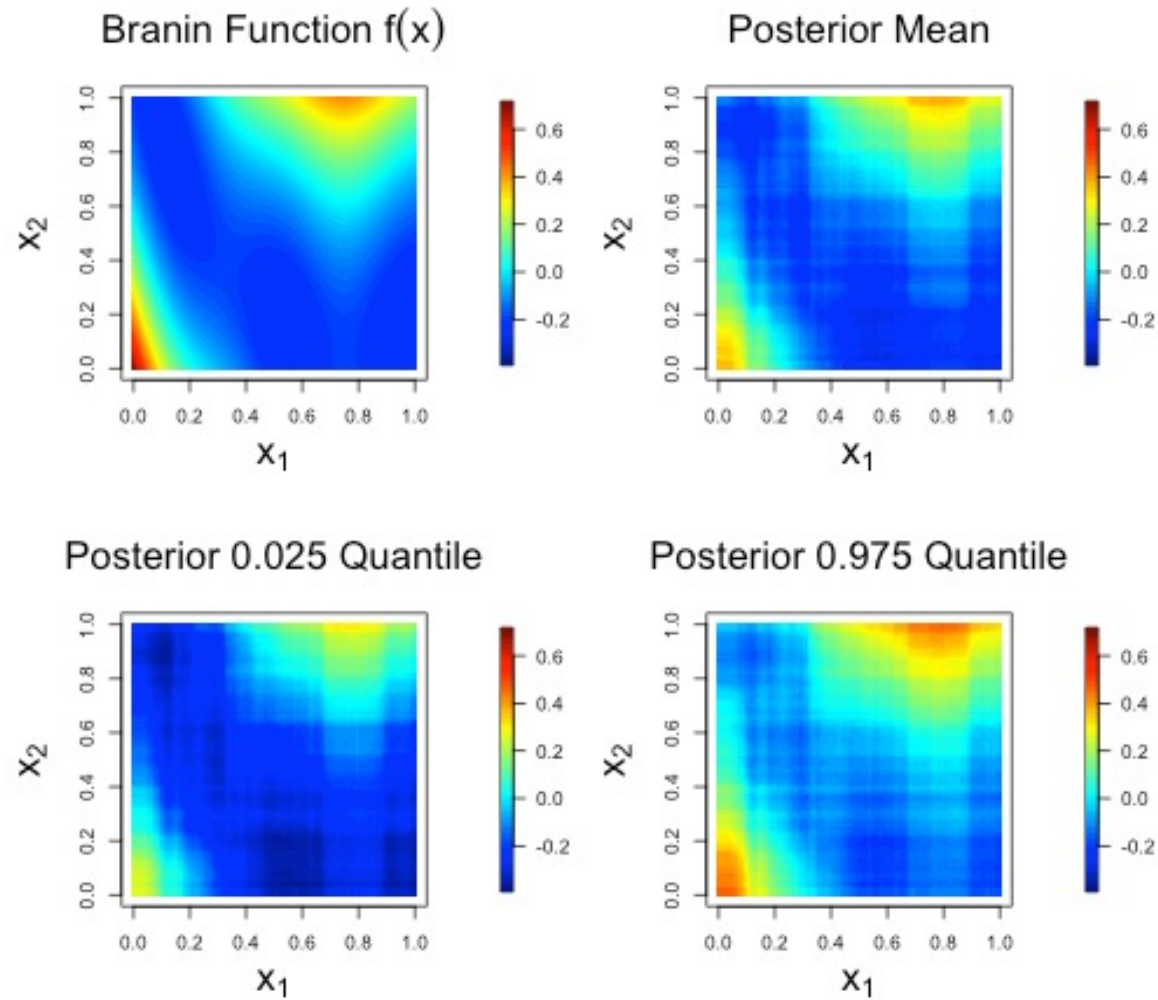
## Posterior Sampling MCMC Algorithm

*[Notation:  $T_{-j} \equiv (T_1, \dots, T_{j-1}, T_{j+1}, \dots, T_m)$  and  $\boldsymbol{\mu}_{-j} \equiv (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_{j-1}, \boldsymbol{\mu}_{j+1}, \dots, \boldsymbol{\mu}_m)$ ]*

For  $i = 1, \dots, N_{mcmc}$ :

1. For  $j = 1, \dots, m$ 
  - a. Sample  $T_j \mid (T_{-j}, \boldsymbol{\mu}_{-j}, \sigma^2, \text{data})$  (Metropolis–Hastings)
  - b. Sample  $\boldsymbol{\mu}_j \mid \cdot$  (Gibbs Step)
2. Sample  $\sigma^2 \mid \cdot$  (Gibbs Step)

# Bayesian Additive Regression Trees



1. Recap of BART
2. Bayesian Inference of the Number of Trees
  - i. Motivation
  - ii. A Fully Bayesian Model
  - iii. Sampling from the Posterior Distribution
  - iv. Code
  - v. Simulations
  - vi. Application to Real Data
3. Conclusion

# How Many Trees???

**Default**  $m = 200$

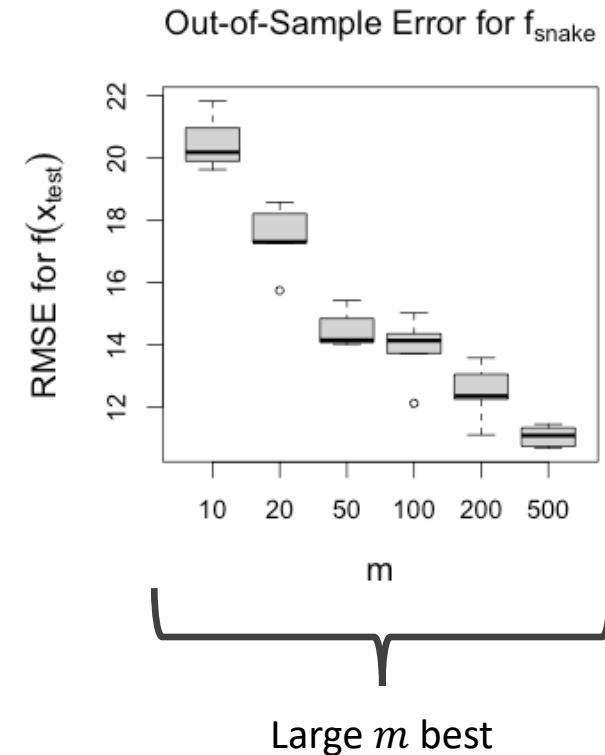
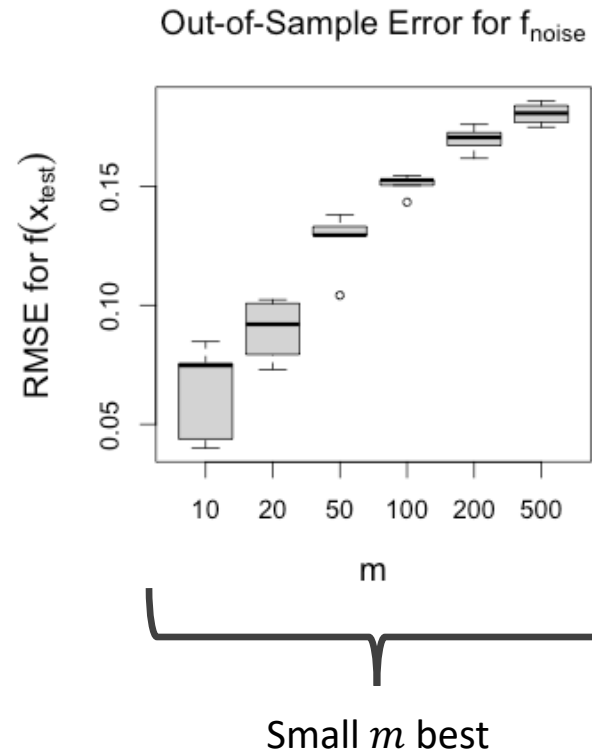
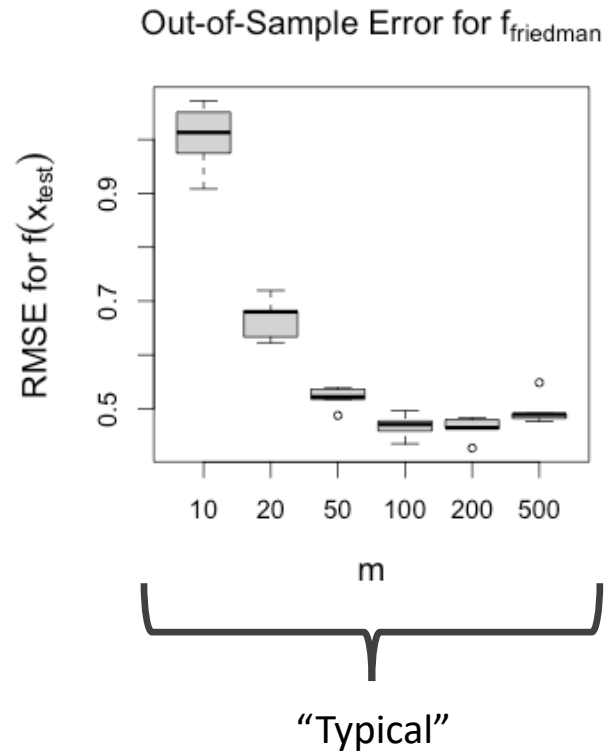
**Large**  $m$

- Flexible Estimation
- More computation
- Risk overfitting

**Small**  $m$

- Improved variable selection
- Less computation
- Risk underfitting

# Out-of-Sample Prediction



$$\text{“RMSE for } f(x_{\text{test}})\text{”} = \sqrt{\frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} (f(x_{\text{test},i}) - \hat{f}(x_{\text{test},i}))^2}$$

# Variable Selection

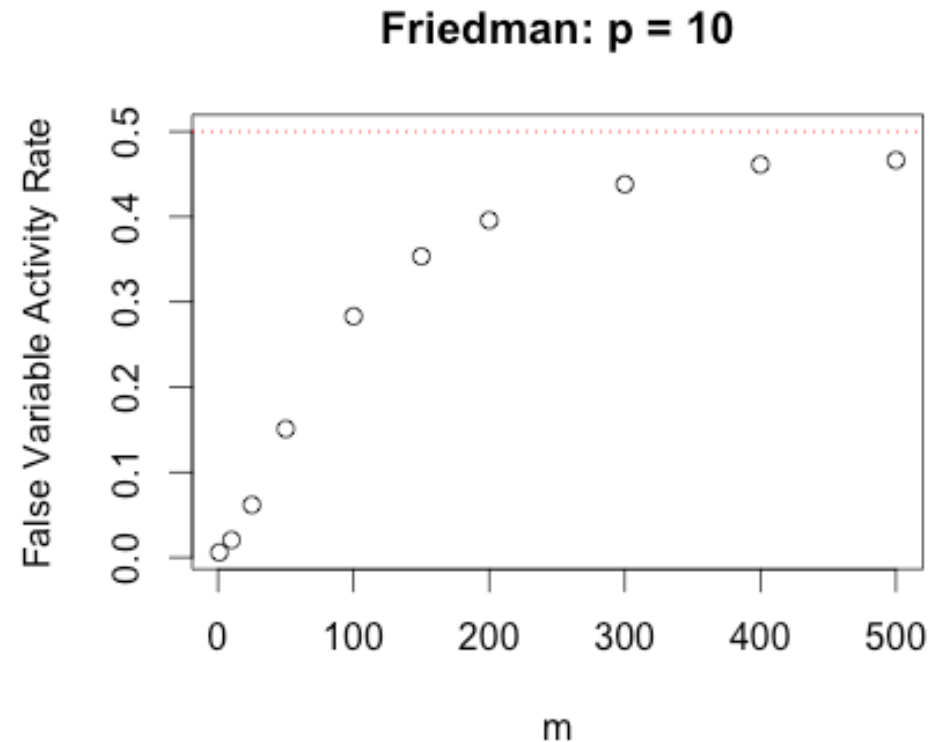
“Real” Inputs



$$f_{\text{friedman}}(\mathbf{x}) = 10 \sin(\pi x_1 x_2) + 20(x_3 - 0.5)^2 + 10x_4 + 5x_5 + \sum_{j=6}^{10} 0x_j$$



“False” Inputs



FVAR = Proportion of branches involving “false” input variables



# Cross-Validation

- Pick a grid of  $m$ -values (e.g.,  $m = 1, 10, 20, 50, 100, 200, 300, 400$ )
- For each value of  $m$ 
  - Split data into train and test sets
  - Fit BART to the training set
  - Predict responses in the test set
- Compare out-of-sample RMSE across the grid
- Pick the value  $m = m_{CV}$  that minimizes RMSE
- Fit a BART model to the full dataset, with  $m = m_{CV}$

# Cross-Validation

## How to pick the grid?

- Pick a grid of  $m$ -values (e.g.,  $m = 1, 10, 20, 50, 100, 200, 300, 400$ )
- For each value of  $m$ 
  - Split data into train and test sets
  - Fit BART to the training set
  - Predict responses in the test set
- Compare out-of-sample RMSE across the grid
- Pick the value  $m = m_{CV}$  that minimizes RMSE
- Fit a BART model to the full dataset, with  $m = m_{CV}$

Expensive!

What about variable selection, computation time, etc.?

1. Recap of BART
- 2. Bayesian Inference of the Number of Trees**
  - i. Motivation
  - ii. A Fully Bayesian Model**
  - iii. Sampling from the Posterior Distribution
  - iv. Code
  - v. Simulations
  - vi. Application to Real Data
3. Conclusion

# Fully Bayesian Inference of $m$

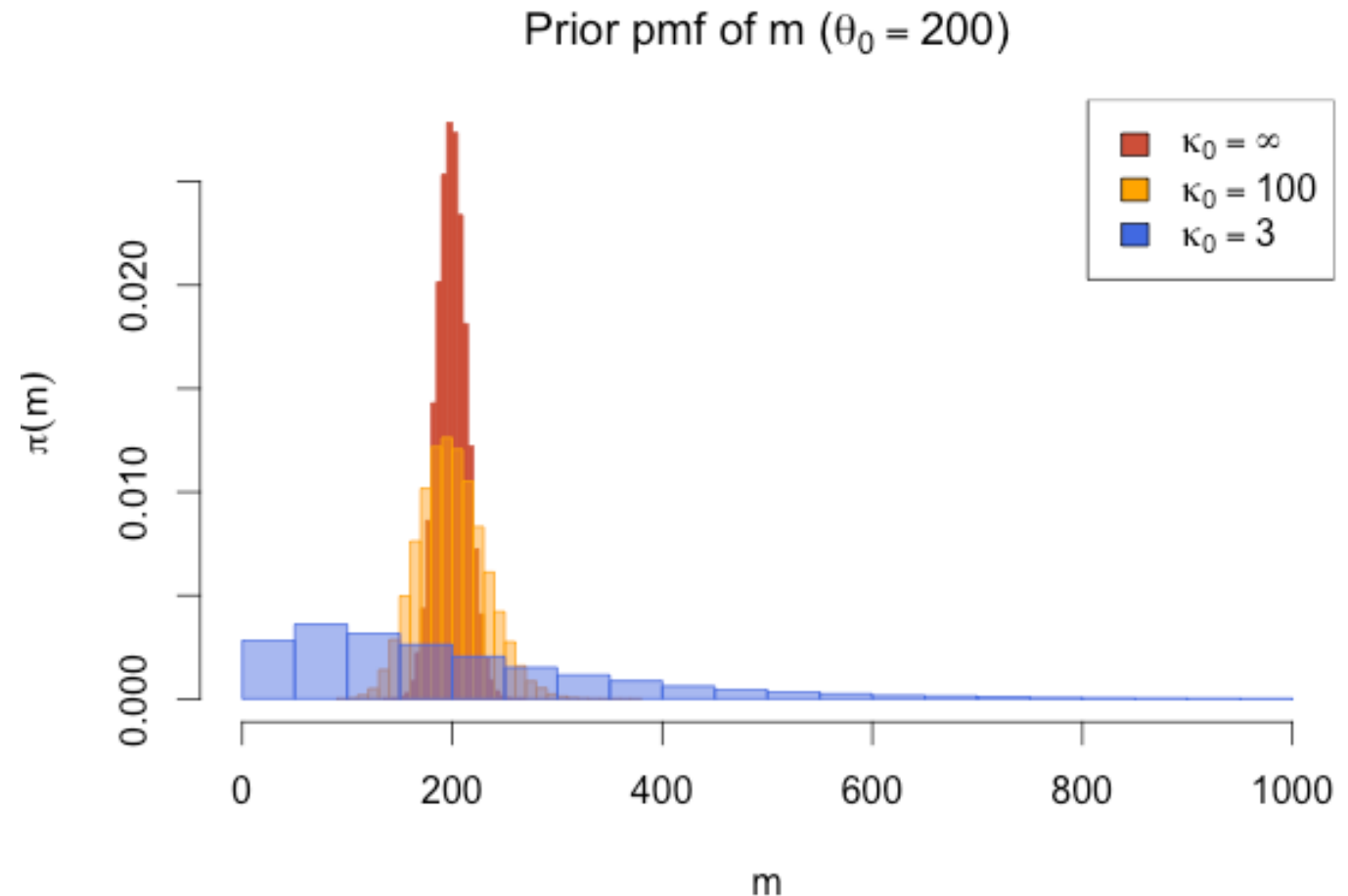
$$\pi(m, T_{1:m}, \boldsymbol{\mu}_{1:m}, \sigma^2) = \pi(\sigma^2)\pi(m) \prod_{j=1}^m \pi(T_j)\pi(\boldsymbol{\mu}_j | T_j, m)$$

- $m \sim \text{Poisson}(\theta)$  [Truncated]
  - Optionally, assign  $\theta$  a hyperprior
- $T_j \sim \text{Tree-Generating Stochastic Process}$  (as before)
- $\mu_{j\ell} | (m, T_j) \sim N(0, \tau_m^2)$ ;  $\ell = 1, \dots, L_j$ ;  $j = 1, \dots, m$  (iid) (as before)
- $\sigma^2 \sim \text{Scaled-inv-}\chi^2(\nu, \lambda)$  (as before)

# Prior Distribution

$$\pi(m) \propto \frac{\theta^m e^{-\theta}}{m!} \mathbf{I}(1 \leq m \leq 1000) \text{ (Truncated Poisson)}$$

- $\Rightarrow E(m) \approx \theta$
- Default  $\theta = 200$
- Optionally, assign  $\theta$  a hyperprior
  - $\theta \sim \frac{\theta_0 \chi_{\kappa_0}^2}{\kappa_0}$
  - $E(\theta) = \theta_0$
  - Degree of Freedom  $\kappa_0$
  - Default  $\theta_0 = 200$



# Prior Distribution

$$\mu_{j\ell} \mid (m, T_j) \sim N(0, \tau_m^2)$$

$$\tau_m = \frac{\max_i y_i - \min_i y_i}{2k\sqrt{m}}$$

$$\Rightarrow f(\mathbf{x}) \sim N\left(0, \left(\frac{\max_i y_i - \min_i y_i}{2k}\right)^2\right) \text{ (for all } m)$$

1. Recap of BART
2. Bayesian Inference of the Number of Trees
  - i. Motivation
  - ii. A Fully Bayesian Model
  - iii. Sampling from the Posterior Distribution
  - iv. Code
  - v. Simulations
  - vi. Application to Real Data
3. Conclusion

# Posterior Sampling MCMC Algorithm

Initialize  $m = m_0$  (default  $m_0 = \theta_0$ )

For  $i = 1, \dots, N_{mcmc}$ :

1. Sample  $\theta \mid m$  (Gibbs Step; if  $\kappa_0 < \infty$ )

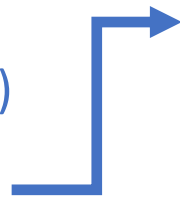


# Posterior Sampling MCMC Algorithm

Initialize  $m = m_0$  (default  $m_0 = \theta_0$ )

For  $i = 1, \dots, N_{mcmc}$ :

1. Sample  $\theta \mid m$  (Gibbs Step; if  $\kappa_0 < \infty$ )
2. Sample  $m \mid \cdot$  (Metropolis-Hastings)



**Randomly select either**

- a) Birth or
- b) Death

# Posterior Sampling MCMC Algorithm

Initialize  $m = m_0$  (default  $m_0 = \theta_0$ )

For  $i = 1, \dots, N_{mcmc}$ :

1. Sample  $\theta \mid m$  (Gibbs Step; if  $\kappa_0 < \infty$ )

2. Sample  $m \mid \cdot$  (Metropolis-Hastings)

- If  $m$  was increased, sample new  $\mu_* \mid \cdot$  (Gibbs Step)

**Randomly select either**

- a) Birth or
- b) Death

# Posterior Sampling MCMC Algorithm

Initialize  $m = m_0$  (default  $m_0 = \theta_0$ )

For  $i = 1, \dots, N_{mcmc}$ :

1. Sample  $\theta \mid m$  (Gibbs Step; if  $\kappa_0 < \infty$ )

2. Sample  $m \mid \cdot$  (Metropolis-Hastings)

- If  $m$  was increased, sample new  $\mu_* \mid \cdot$  (Gibbs Step)

3. For  $j = 1, \dots, m$

a. Sample  $T_j \mid (T_{-j}, \mu_{-j}, \sigma^2, \text{data})$  (Metropolis-Hastings)

b. Sample  $\mu_j \mid \cdot$  (Gibbs Step)

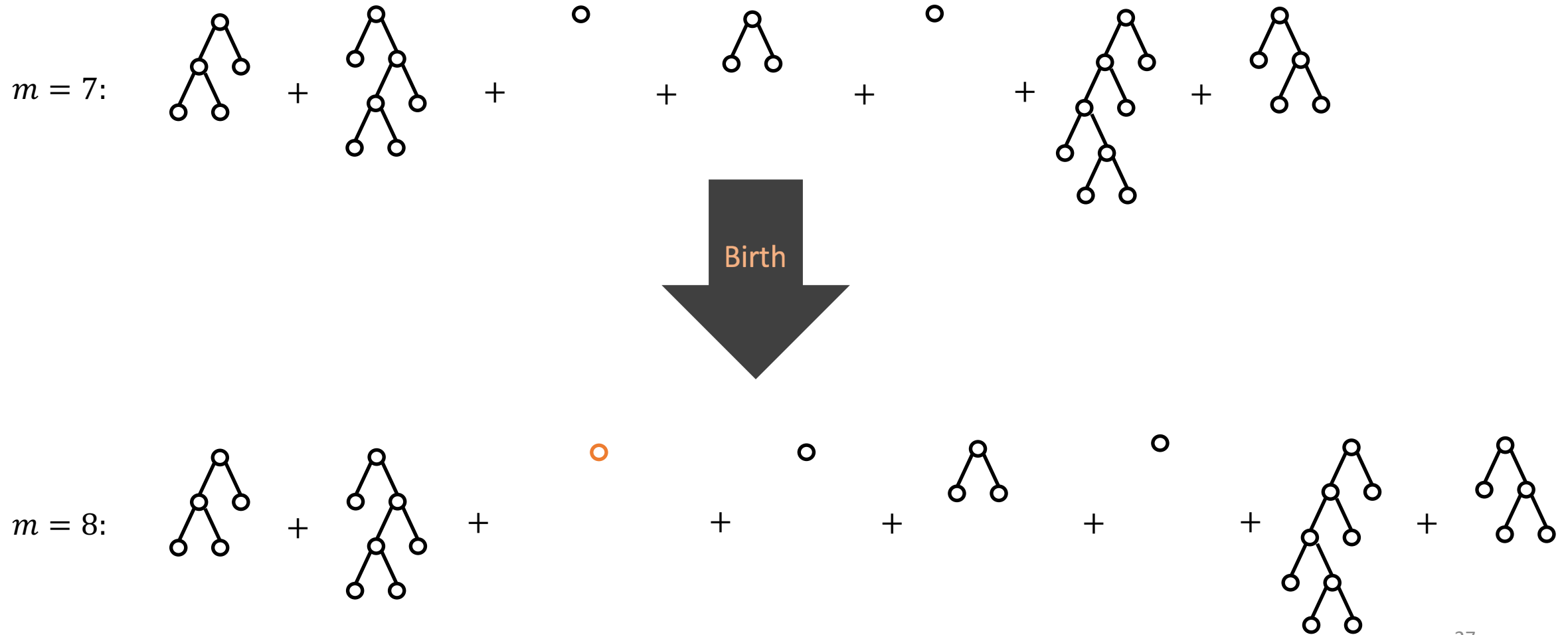
4. Sample  $\sigma^2 \mid \cdot$  (Gibbs Step)

Randomly select either

- a) Birth or
- b) Death

Same as  
Standard BART

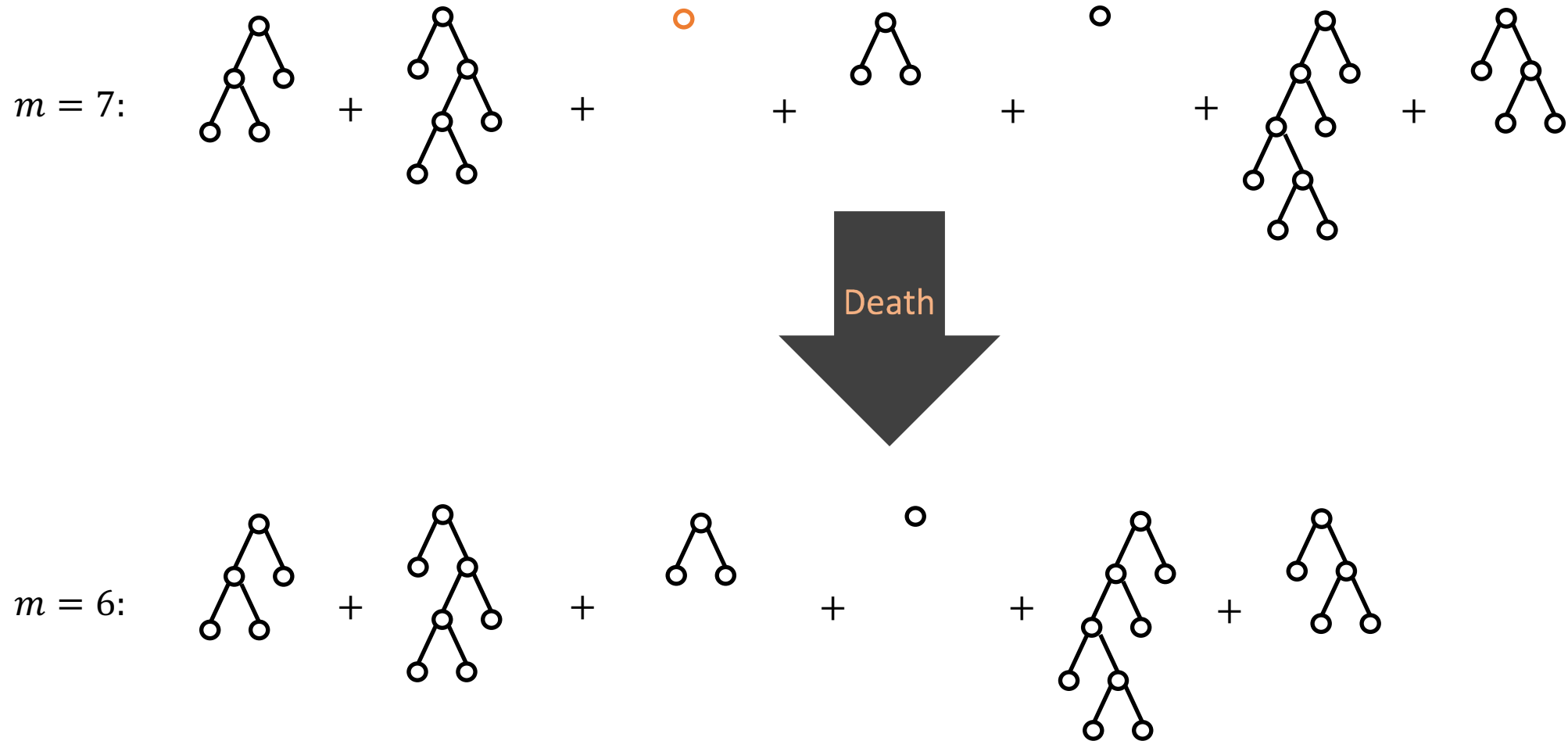
# Birth Transition



# Birth Transition

- Why stumps?
- Why randomize the location of the new tree?
  - Trees are exchangeable, but ordered in the prior distribution
  - Need reversibility

# Death Transition



# MH Transition

- Current parameters  $\psi = (m, T_{1:m}, \boldsymbol{\mu}_{1:m}, \sigma^2)$
- Randomly select either birth or death transition ( $\Pr(\text{birth}) = \Pr(\text{death}) = 0.5$ )

- **Birth Transition**

- Select a location to insert stump  $T^*$

$$q(\psi \rightarrow \psi^{\text{birth}}) = 0.5 \times \frac{1}{m + 1}$$

- Update to  $\psi \rightarrow \psi^{\text{birth}} = (m + 1, T_{1:(m+1)}^*, \boldsymbol{\mu}_{1:m}, \sigma^2)$

- **Death Transition (if there are any “stumps”)**

- Select a stump  $T^*$  to delete

$$q(\psi \rightarrow \psi^{\text{death}}) = 0.5 \times \frac{1}{m_{\text{stumps}}}$$

- Update to  $\psi \rightarrow \psi^{\text{death}} = (m - 1, T_{1:(m-1)}^*, \boldsymbol{\mu}_{1:(m-1)}, \sigma^2)$

- Accept with the MH acceptance probability:  $\min\{1, \text{MH Ratio}\}$

# MH Acceptance Probability: $\min\{1, \text{MH Ratio}\}$

## MH Ratio for Birth

$$= \frac{\pi(\psi^{\text{birth}} \mid \text{data})}{\pi(\psi \mid \text{data})} \times \frac{q(\psi^{\text{birth}} \rightarrow \psi)}{q(\psi \rightarrow \psi^{\text{birth}})}$$

Ratio of Posterior Distributions

Ratio of Transition Probabilities



# MH Acceptance Probability: $\min\{1, \text{MH Ratio}\}$

## MH Ratio for Birth

$$= \frac{\pi(\psi^{\text{birth}} \mid \text{data})}{\pi(\psi \mid \text{data})} \times \frac{q(\psi^{\text{birth}} \rightarrow \psi)}{q(\psi \rightarrow \psi^{\text{birth}})}$$

$$= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \underbrace{\frac{\pi(\sigma^2)}{\pi(\sigma^2)} \times \frac{\pi(m+1)}{\pi(m)} \times \frac{\pi(T^*) \prod_j^m \pi(T_j)}{\prod_j^m \pi(T_j)} \times \frac{\prod_j^m \pi(\mu_j \mid m+1, T_j)}{\prod_j^m \pi(\mu_j \mid m, T_j)}}_{\text{Prior Distribution}} \times \frac{0.5/(m_{\text{stumps}} + 1)}{0.5/(m + 1)}$$

Likelihood

Prior Distribution

Ratio of Transition Probabilities

# MH Acceptance Probability: $\min\{1, \text{MH Ratio}\}$

## MH Ratio for Birth

$$\begin{aligned} &= \frac{\pi(\psi^{\text{birth}} \mid \text{data})}{\pi(\psi \mid \text{data})} \times \frac{q(\psi^{\text{birth}} \rightarrow \psi)}{q(\psi \rightarrow \psi^{\text{birth}})} \\ &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \frac{\pi(\sigma^2)}{\pi(\sigma^2)} \times \frac{\pi(m+1)}{\pi(m)} \times \frac{\pi(T^*) \prod_j^m \pi(T_j)}{\prod_j^m \pi(T_j)} \times \frac{\prod_j^m \pi(\mu_j \mid m+1, T_j)}{\prod_j^m \pi(\mu_j \mid m, T_j)} \times \frac{0.5/(m_{\text{stumps}} + 1)}{0.5/(m+1)} \\ &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \frac{\pi(m+1)}{\pi(m)} \times \pi(T^*) \times \frac{\prod_j^m \pi(\mu_j \mid m+1, T_j)}{\prod_j^m \pi(\mu_j \mid m, T_j)} \times \frac{1/(m_{\text{stumps}} + 1)}{1/(m+1)} \end{aligned}$$

# MH Acceptance Probability: $\min\{1, \text{MH Ratio}\}$

## MH Ratio for Birth

$$\begin{aligned} &= \frac{\pi(\psi^{\text{birth}} \mid \text{data})}{\pi(\psi \mid \text{data})} \times \frac{q(\psi^{\text{birth}} \rightarrow \psi)}{q(\psi \rightarrow \psi^{\text{birth}})} \\ &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \frac{\pi(\sigma^2)}{\pi(\sigma^2)} \times \frac{\pi(m+1)}{\pi(m)} \times \frac{\pi(T^*) \prod_j^m \pi(T_j)}{\prod_j^m \pi(T_j)} \times \frac{\prod_j^m \pi(\mu_j \mid m+1, T_j)}{\prod_j^m \pi(\mu_j \mid m, T_j)} \times \frac{0.5/(m_{\text{stumps}}+1)}{0.5/(m+1)} \\ &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \frac{\pi(m+1)}{\pi(m)} \times \pi(T^*) \times \frac{\prod_j^m \pi(\mu_j \mid m+1, T_j)}{\prod_j^m \pi(\mu_j \mid m, T_j)} \times \frac{1/(m_{\text{stumps}}+1)}{1/(m+1)} \\ &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \frac{\theta}{m+1} \times \end{aligned}$$

# MH Acceptance Probability: $\min\{1, \text{MH Ratio}\}$

## MH Ratio for Birth

$$\begin{aligned} &= \frac{\pi(\psi^{\text{birth}} \mid \text{data})}{\pi(\psi \mid \text{data})} \times \frac{q(\psi^{\text{birth}} \rightarrow \psi)}{q(\psi \rightarrow \psi^{\text{birth}})} \\ &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \frac{\pi(\sigma^2)}{\pi(\sigma^2)} \times \frac{\pi(m+1)}{\pi(m)} \times \frac{\pi(T^*) \prod_j^m \pi(T_j)}{\prod_j^m \pi(T_j)} \times \frac{\prod_j^m \pi(\mu_j \mid m+1, T_j)}{\prod_j^m \pi(\mu_j \mid m, T_j)} \times \frac{0.5/(m_{\text{stumps}}+1)}{0.5/(m+1)} \\ &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \frac{\pi(m+1)}{\pi(m)} \times \pi(T^*) \times \frac{\prod_j^m \pi(\mu_j \mid m+1, T_j)}{\prod_j^m \pi(\mu_j \mid m, T_j)} \times \frac{1/(m_{\text{stumps}}+1)}{1/(m+1)} \\ &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \frac{\theta}{m+1} \times \pi(T^*) \times \end{aligned}$$

# MH Acceptance Probability: $\min\{1, \text{MH Ratio}\}$

## MH Ratio for Birth

$$\begin{aligned} &= \frac{\pi(\psi^{\text{birth}} \mid \text{data})}{\pi(\psi \mid \text{data})} \times \frac{q(\psi^{\text{birth}} \rightarrow \psi)}{q(\psi \rightarrow \psi^{\text{birth}})} \\ &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \frac{\pi(\sigma^2)}{\pi(\sigma^2)} \times \frac{\pi(m+1)}{\pi(m)} \times \frac{\pi(T^*) \prod_j^m \pi(T_j)}{\prod_j^m \pi(T_j)} \times \frac{\prod_j^m \pi(\mu_j \mid m+1, T_j)}{\prod_j^m \pi(\mu_j \mid m, T_j)} \times \frac{0.5/(m_{\text{stumps}}+1)}{0.5/(m+1)} \\ &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \frac{\pi(m+1)}{\pi(m)} \times \pi(T^*) \times \frac{\prod_j^m \pi(\mu_j \mid m+1, T_j)}{\prod_j^m \pi(\mu_j \mid m, T_j)} \times \frac{1/(m_{\text{stumps}}+1)}{1/(m+1)} \\ &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \frac{\theta}{m+1} \times \pi(T^*) \times \left(\frac{m}{m+1}\right)^{-\sum_{j=1}^m \frac{L_j}{2}} \exp\left(-\frac{1}{2m\tau_m^2} \sum_j^m \sum_{\ell=1}^{L_j} \mu_{j\ell}^2\right) \end{aligned}$$

# MH Acceptance Probability: $\min\{1, \text{MH Ratio}\}$

## MH Ratio for Birth

$$\begin{aligned} &= \frac{\pi(\psi^{\text{birth}} \mid \text{data})}{\pi(\psi \mid \text{data})} \times \frac{q(\psi^{\text{birth}} \rightarrow \psi)}{q(\psi \rightarrow \psi^{\text{birth}})} \\ &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \frac{\pi(\sigma^2)}{\pi(\sigma^2)} \times \frac{\pi(m+1)}{\pi(m)} \times \frac{\pi(T^*) \prod_j^m \pi(T_j)}{\prod_j^m \pi(T_j)} \times \frac{\prod_j^m \pi(\mu_j \mid m+1, T_j)}{\prod_j^m \pi(\mu_j \mid m, T_j)} \times \frac{0.5/(m_{\text{stumps}}+1)}{0.5/(m+1)} \\ &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \frac{\pi(m+1)}{\pi(m)} \times \pi(T^*) \times \frac{\prod_j^m \pi(\mu_j \mid m+1, T_j)}{\prod_j^m \pi(\mu_j \mid m, T_j)} \times \frac{1/(m_{\text{stumps}}+1)}{1/(m+1)} \\ &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \times \frac{\theta}{m+1} \times \pi(T^*) \times \left(\frac{m}{m+1}\right)^{-\sum_{j=1}^m \frac{L_j}{2}} \exp\left(-\frac{1}{2m\tau_m^2} \sum_j^m \sum_{\ell=1}^{L_j} \mu_{j\ell}^2\right) \times \frac{m+1}{m_{\text{stumps}}+1} \end{aligned}$$

# Likelihood Ratio

$$\text{Likelihood Ratio} = \frac{\pi(\text{data} | \psi^{\text{birth}})}{\pi(\text{data} | \psi)}$$

# Likelihood Ratio

$$\text{Likelihood Ratio} = \frac{\pi(\text{data} | \psi^{\text{birth}})}{\pi(\text{data} | \psi)}$$

$$= \frac{\pi(\text{data} | m + 1, T_{1:(m+1)}^*, \boldsymbol{\mu}_{1:m}, \sigma^2)}{\pi(\text{data} | m, T_{1:m}, \boldsymbol{\mu}_{1:m}, \sigma^2)}$$

Marginal Likelihood (m+1 trees)

Full likelihood (m trees)



# Likelihood Ratio

$$\text{Likelihood Ratio} = \frac{\pi(\text{data} | \psi^{\text{birth}})}{\pi(\text{data} | \psi)}$$

$$= \frac{\pi(\text{data} | m + 1, T_{1:(m+1)}^*, \boldsymbol{\mu}_{1:m}, \sigma^2)}{\pi(\text{data} | m, T_{1:m}, \boldsymbol{\mu}_{1:m}, \sigma^2)}$$

$$= \frac{\int_{\mathbb{R}} \pi(\text{data} | m + 1, T_{1:(m+1)}^*, \boldsymbol{\mu}_{1:m}, \mu^*, \sigma^2) \pi(\mu^* | m + 1, T^*) d\mu^*}{\pi(\text{data} | m, T_{1:m}, \boldsymbol{\mu}_{1:m}, \sigma^2)}$$

Full Likelihood (m+1)



Prior Distribution of  $\mu^*$



# Likelihood Ratio

$$\begin{aligned}\text{Likelihood Ratio} &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \\ &= \frac{\pi(\text{data} \mid m + 1, T_{1:(m+1)}^*, \boldsymbol{\mu}_{1:m}, \sigma^2)}{\pi(\text{data} \mid m, T_{1:m}, \boldsymbol{\mu}_{1:m}, \sigma^2)} \\ &= \frac{\int_{\mathbb{R}} \pi(\text{data} \mid m + 1, T_{1:(m+1)}^*, \boldsymbol{\mu}_{1:m}, \mu^*, \sigma^2) \pi(\mu^* \mid m, T^*) d\mu^*}{\pi(\text{data} \mid m, T_{1:m}, \boldsymbol{\mu}_{1:m}, \sigma^2)} \\ &= \frac{\int_{\mathbb{R}} \prod_{i=1}^n N(y_i; \mu^* + \sum_{j=1}^m g(\mathbf{x}_i; T_j, \boldsymbol{\mu}_j), \sigma^2) N(\mu^*; 0, \tau_m^2) d\mu^*}{\prod_{i=1}^n N(y_i; \sum_{j=1}^m g(\mathbf{x}_i; T_j, \boldsymbol{\mu}_j), \sigma^2)}\end{aligned}$$

# Likelihood Ratio

$$\begin{aligned}\text{Likelihood Ratio} &= \frac{\pi(\text{data} \mid \psi^{\text{birth}})}{\pi(\text{data} \mid \psi)} \\ &= \frac{\pi(\text{data} \mid m+1, T_{1:(m+1)}^*, \boldsymbol{\mu}_{1:m}, \sigma^2)}{\pi(\text{data} \mid m, T_{1:m}, \boldsymbol{\mu}_{1:m}, \sigma^2)} \\ &= \frac{\int_{\mathbb{R}} \pi(\text{data} \mid m+1, T_{1:(m+1)}^*, \boldsymbol{\mu}_{1:m}, \mu^*, \sigma^2) \pi(\mu^* \mid m, T^*) d\mu^*}{\pi(\text{data} \mid m, T_{1:m}, \boldsymbol{\mu}_{1:m}, \sigma^2)} \\ &= \frac{\int_{\mathbb{R}} \prod_{i=1}^n N(y_i; \mu^* + \sum_{j=1}^m g(\mathbf{x}_i; T_j, \boldsymbol{\mu}_j), \sigma^2) N(\mu^*; 0, \tau_m^2) d\mu^*}{\prod_{i=1}^n N(y_i; \sum_{j=1}^m g(\mathbf{x}_i; T_j, \boldsymbol{\mu}_j), \sigma^2)} \\ &= \left( \frac{\sigma^2}{n\tau_{m+1}^2 + \sigma^2} + 1 \right)^{1/2} \exp \left( \frac{n^2 \tau_{m+1}^2 (\sum_{i=1}^n [y_i - \sum_{j=1}^m g(\mathbf{x}_i; T_j, \boldsymbol{\mu}_j)])^2}{2\sigma^2 (n\tau_{m+1}^2 + \sigma^2)} \right)\end{aligned}$$

1. Recap of BART
- 2. Bayesian Inference of the Number of Trees**
  - i. Motivation
  - ii. A Fully Bayesian Model
  - iii. Sampling from the Posterior Distribution
  - iv. Code**
  - v. Simulations
  - vi. Application to Real Data
3. Conclusion

# Code

- R implementation forthcoming:
  - `bart(X, y, learntree = TRUE, ntreemean = 200, ntreedf = Inf)`

$$m \sim \text{Pois}(\theta) I(1 \leq m \leq 1000)$$

$$\theta \sim \frac{\theta_0 \chi_{\kappa_0}^2}{\kappa_0}$$

1. Recap of BART
- 2. Bayesian Inference of the Number of Trees**
  - i. Motivation
  - ii. A Fully Bayesian Model
  - iii. Sampling from the Posterior Distribution
  - iv. Code
  - v. Simulations**
  - vi. Application to Real Data
3. Conclusion

# Simulation Setup

## Fully Bayesian Inference for $m$

- Generate training and test data
  - $\mathbf{x}_i \sim \text{Unif}(0,1)^p$  and
  - $y_i | \mathbf{x}_i \sim N(f(\mathbf{x}_i), 1)$
- For  $\kappa_0 \in \{3, 100, \infty\}$  (with  $\theta_0 = 200$ )
  - Fit BART to training set with Bayesian inference for  $m$

## Cross-Validation

- For  $m \in \text{Grid}$ :
  - Generate training and test sets
    - $\mathbf{x}_i \sim \text{Unif}(0,1)^p$  and
    - $y_i | \mathbf{x}_i \sim N(f(\mathbf{x}_i), 1)$
  - Fit BART to training set with  $m$  trees


Compare accuracy using “RMSE for  $f(\mathbf{x}_{\text{test}})$ ” = 
$$\sqrt{\frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} (f(\mathbf{x}_{\text{test},i}) - \hat{f}(\mathbf{x}_{\text{test},i}))^2}$$

# Simulation Setup

Bayesian Inference

Cross-Validation

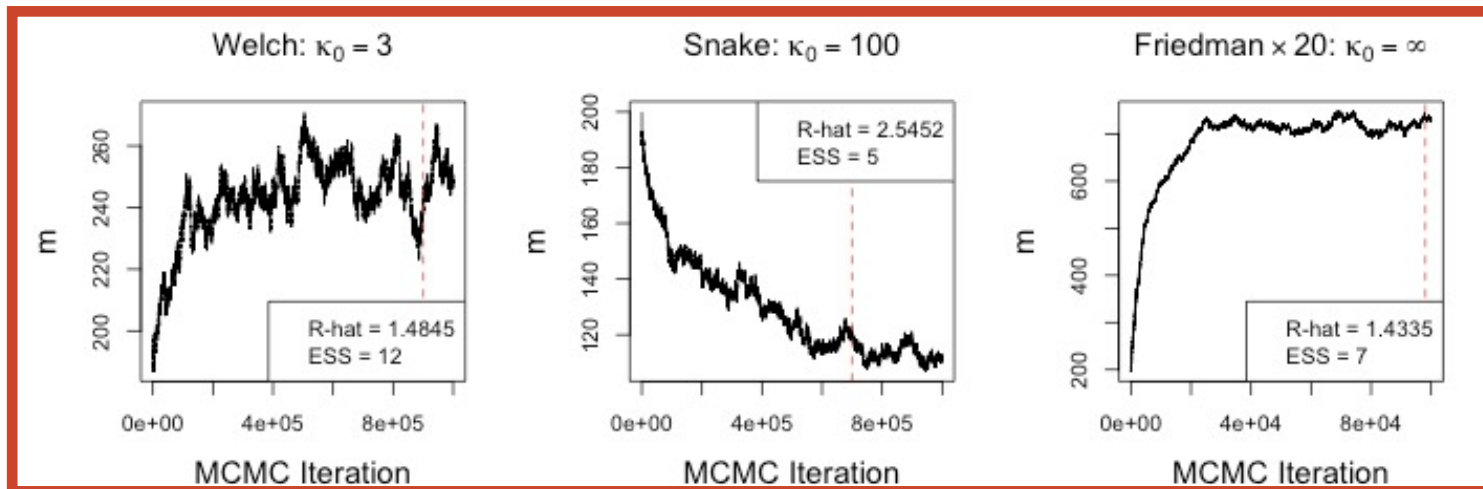
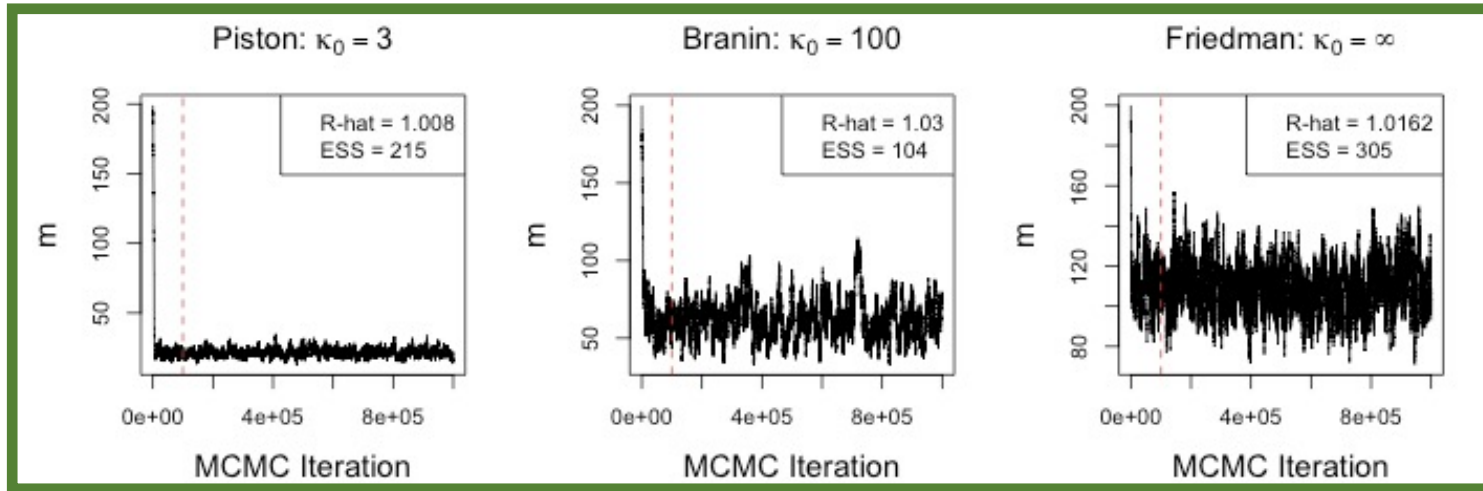
## Simulations



$f$	$n_{\text{train}}$	$n_{\text{test}}$	$p$	SNR	$N_{\text{mc}}$ (infer $m$ )	$N_{\text{mc}}$ (fix $m$ )
Friedman	500	1,000	10	23.8	1,000,000 (✓)	3,000
Borehole	500	1,000	8	20.9	1,000,000 (✓)	3,000
Branin	1,000	2,000	2	18.2	1,000,000 (✓)	3,000
Piston	1,500	3,000	7	20.6	1,000,000 (✓)	3,000
Snake	10,000	10,000	2	2930	1,000,000 (✗)	100,000
Welch	10,000	10,000	20	2809	1,000,000 (✗)	100,000
Friedman×20	10,000	10,000	100	23.8	100,000 (✗)	50,000
300-Step	12,000	12,000	300	100	200,000 (✗)	10,000
100-Step	4,000	8,000	100	100	1,000,000 (✗)	10,000
1-Step	100	200	1	100	1,000,000 (✓)	3,000
$T_4$	800	1,000	15	340	1,000,000 (✗)	30,000

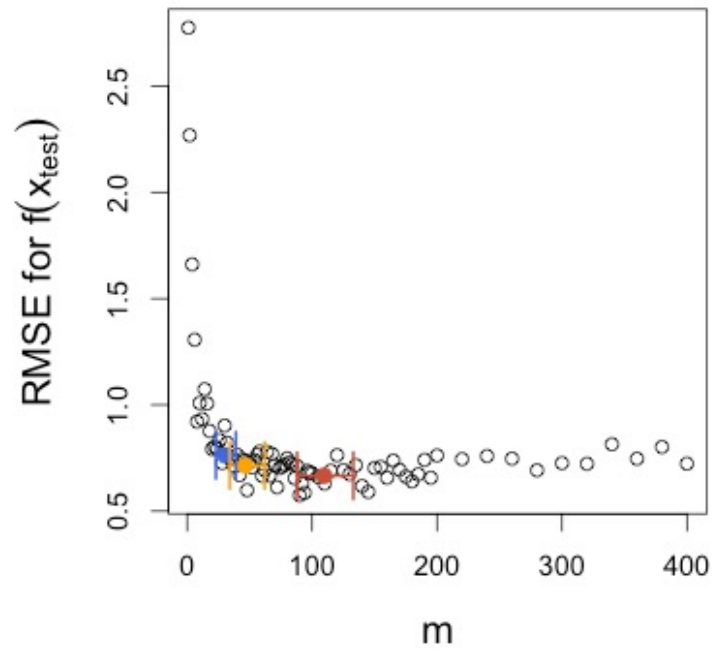


# Convergence of $m$

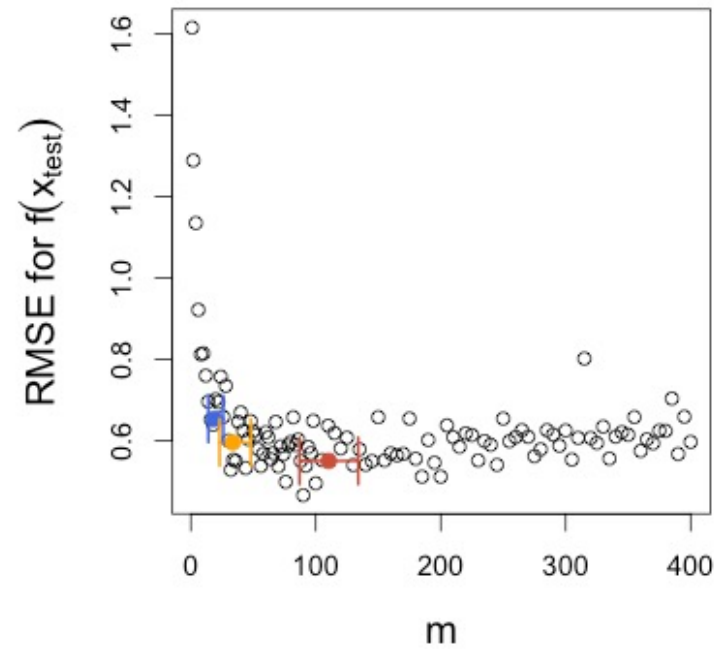


# Results

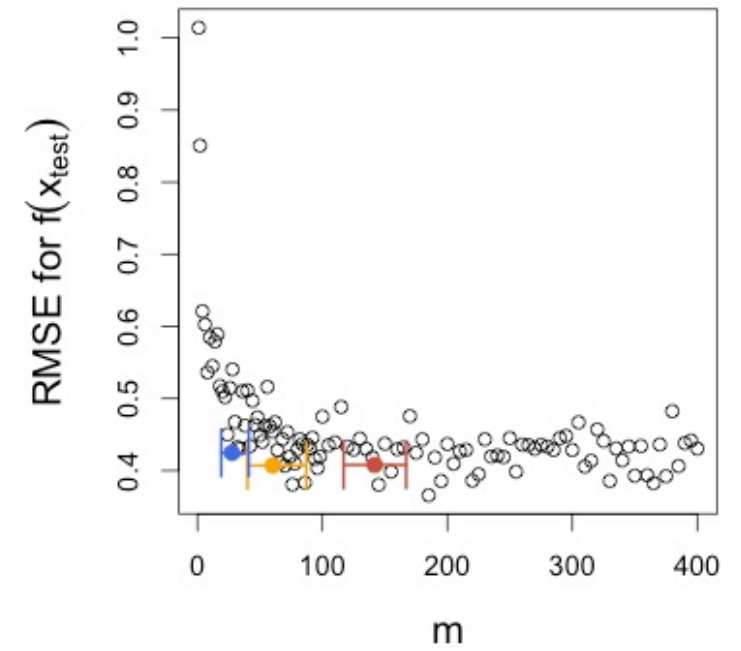
Friedman



Borehole



Branin



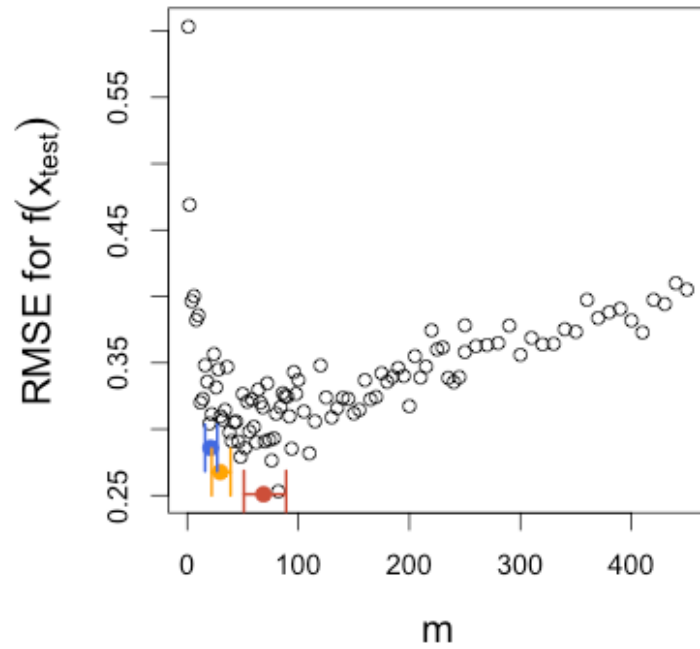
- $\kappa_0 = 3$
- $\kappa_0 = 100$
- $\kappa_0 = \infty$
- Fixed  $m$

# Results

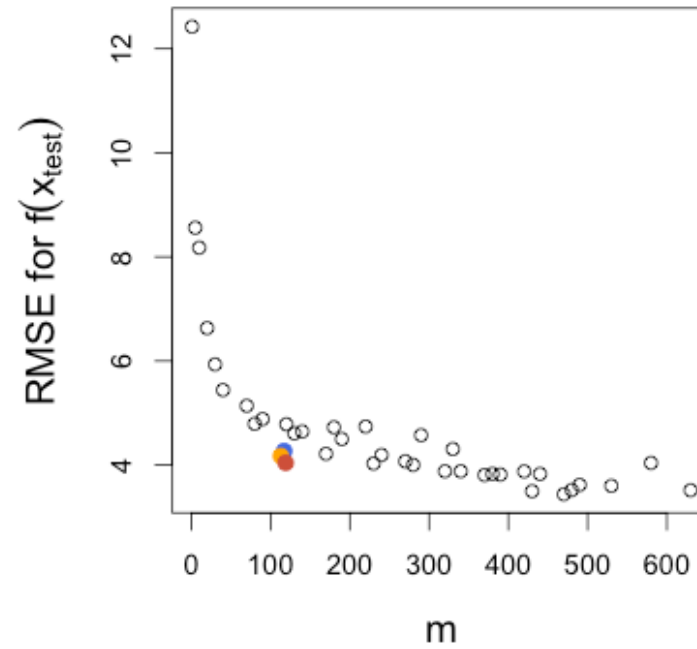
$p = 100$  inputs: all active!  
 $f$  has an additive structure



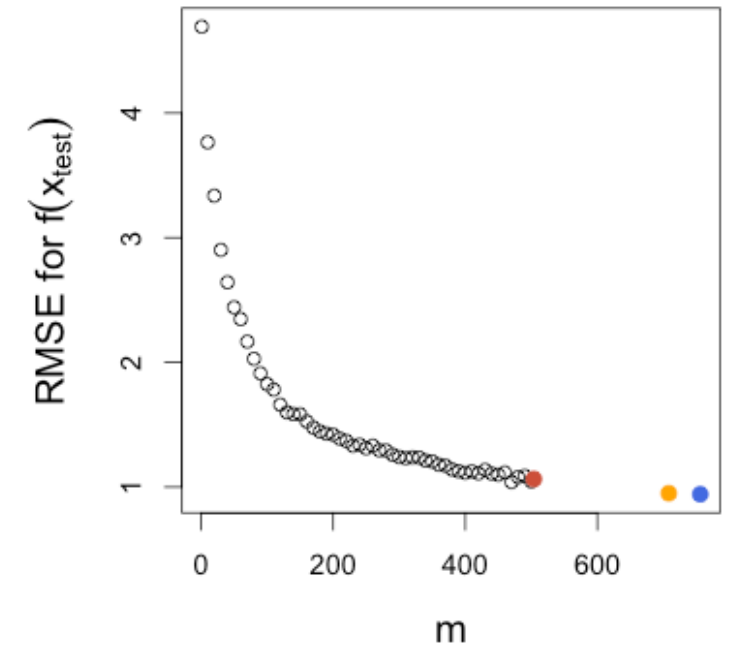
Piston



Snake



Friedman  $\times 20$



- $\kappa_0 = 3$
- $\kappa_0 = 100$
- $\kappa_0 = \infty$
- Fixed  $m$

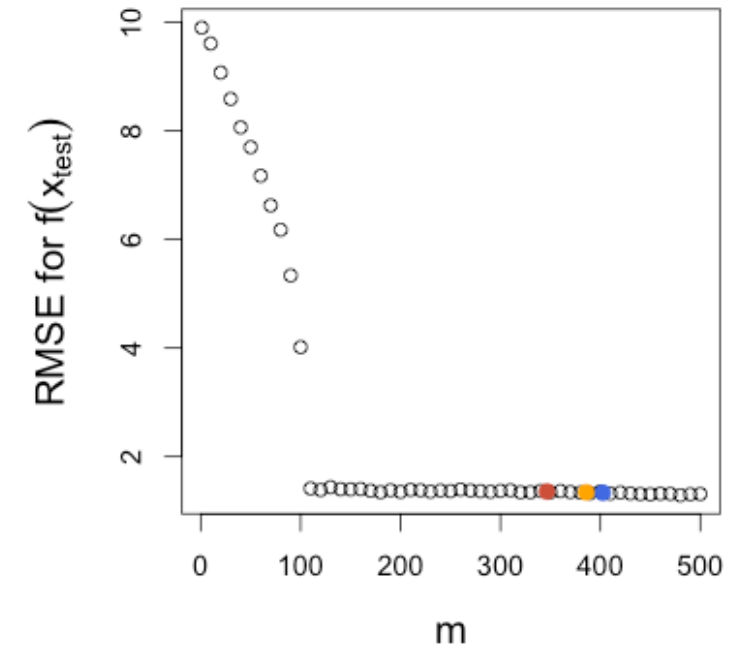
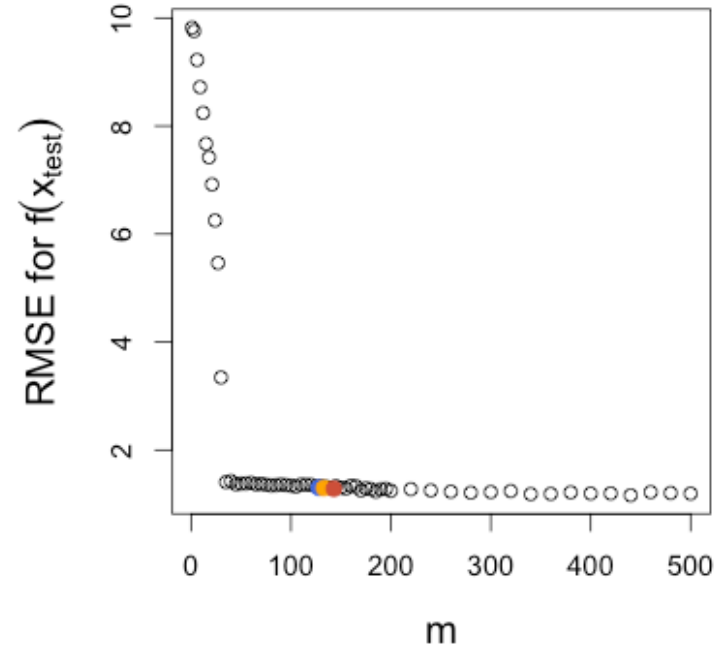
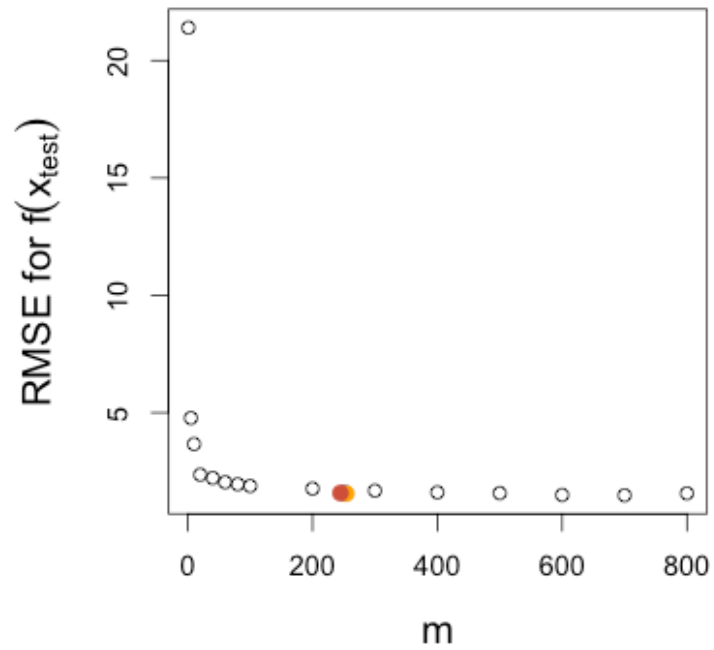
# Results

$$p\text{-step: } f(\mathbf{x}) = \frac{20}{\sqrt{p}} \sum_{j=1}^p \mathbf{I}(\mathbf{x}_j \geq 0.5)$$

100-Step

300-Step

Welch



- $\kappa_0 = 3$
- $\kappa_0 = 100$
- $\kappa_0 = \infty$
- Fixed  $m$

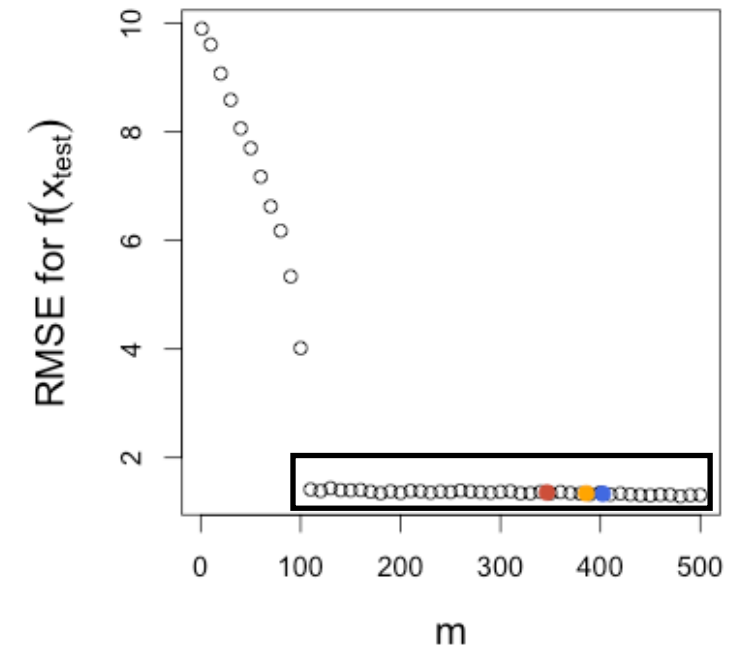
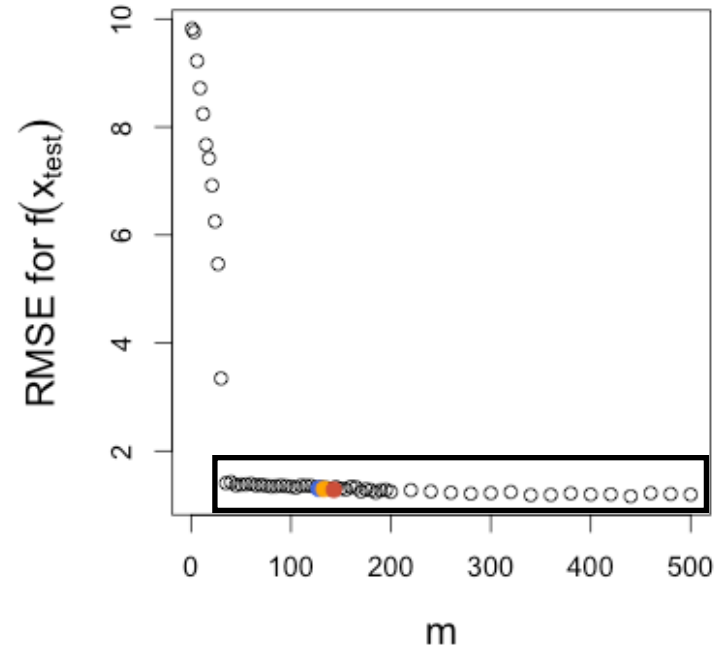
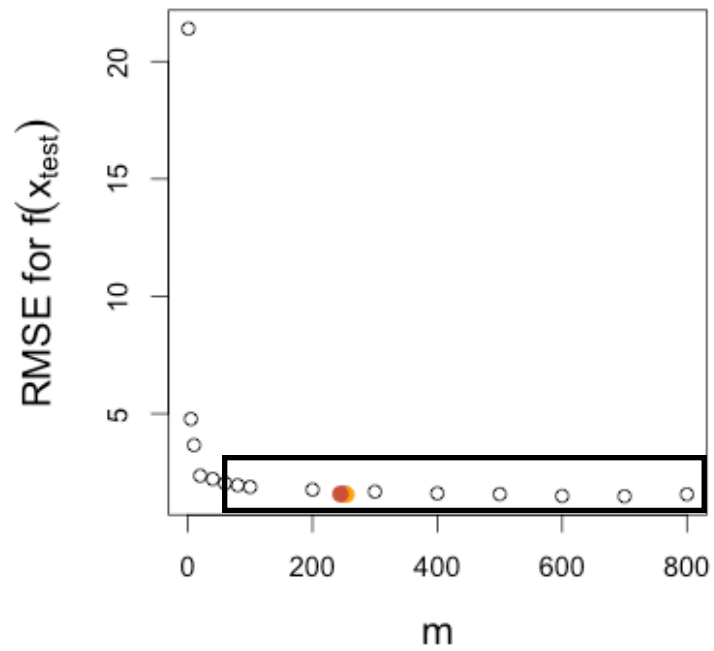
# Results

$$p\text{-step: } f(\mathbf{x}) = \frac{20}{\sqrt{p}} \sum_{j=1}^p \mathbf{I}(\mathbf{x}_j \geq 0.5)$$

100-Step

300-Step

Welch



- $\kappa_0 = 3$
- $\kappa_0 = 100$
- $\kappa_0 = \infty$
- Fixed  $m$

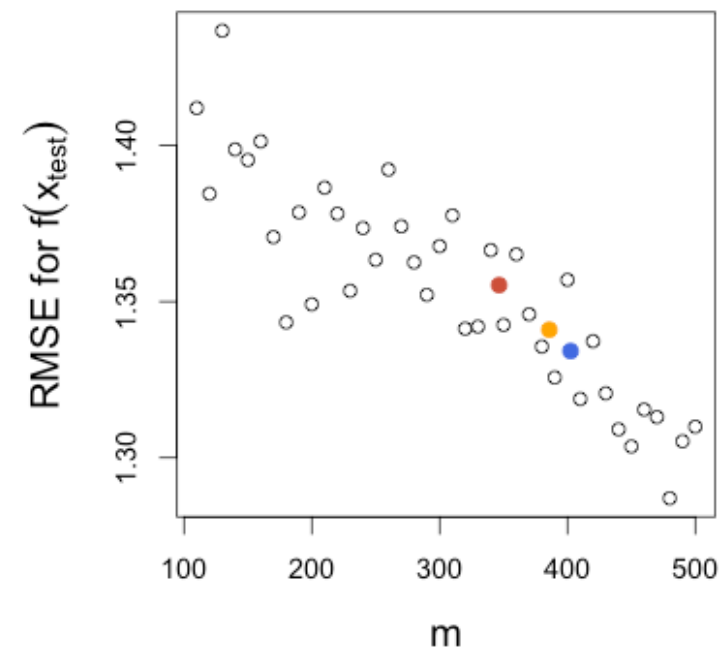
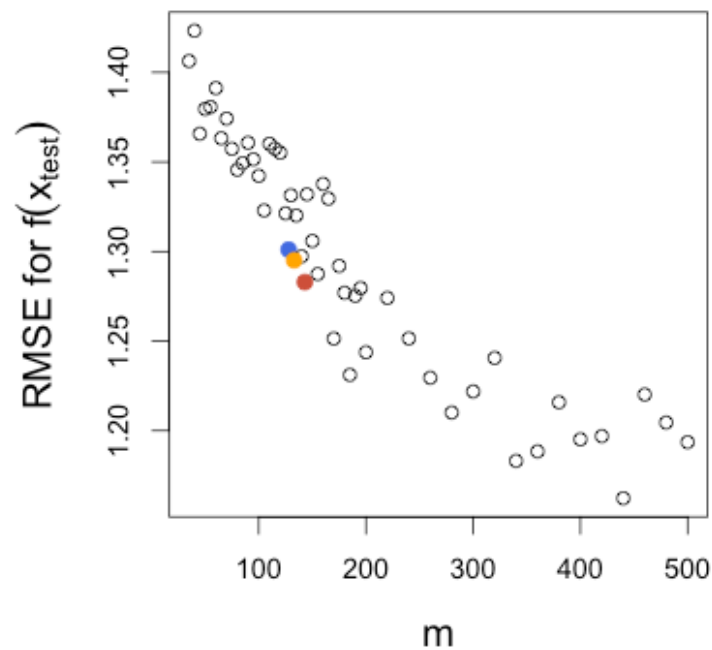
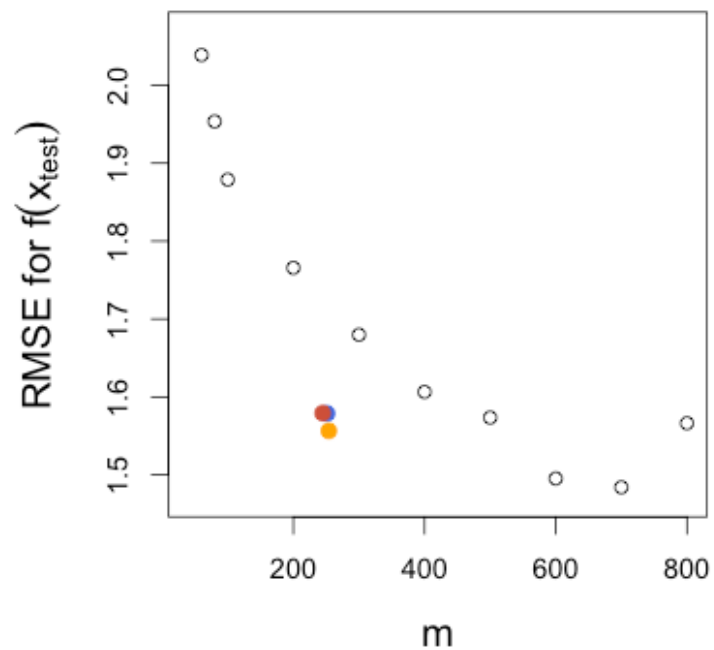
# Results

$$p\text{-step: } f(\mathbf{x}) = \frac{20}{\sqrt{p}} \sum_{j=1}^p \mathbf{I}(\mathbf{x}_j \geq 0.5)$$

100-Step

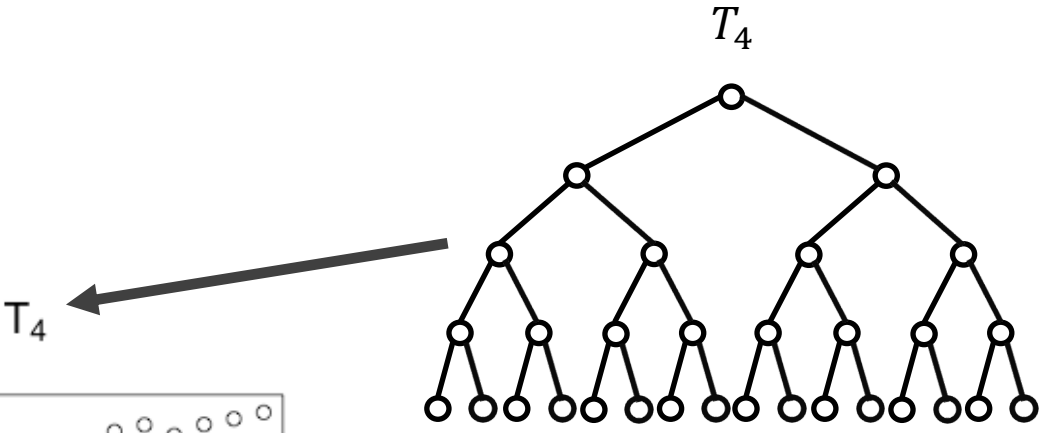
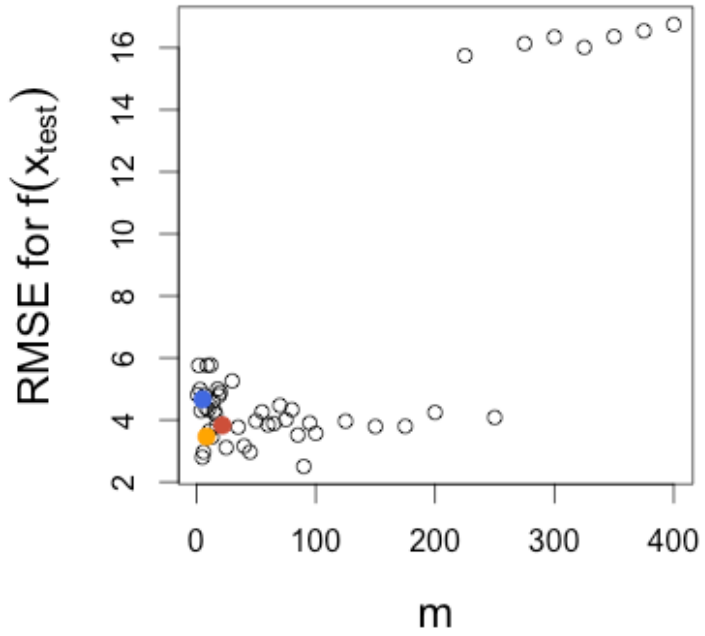
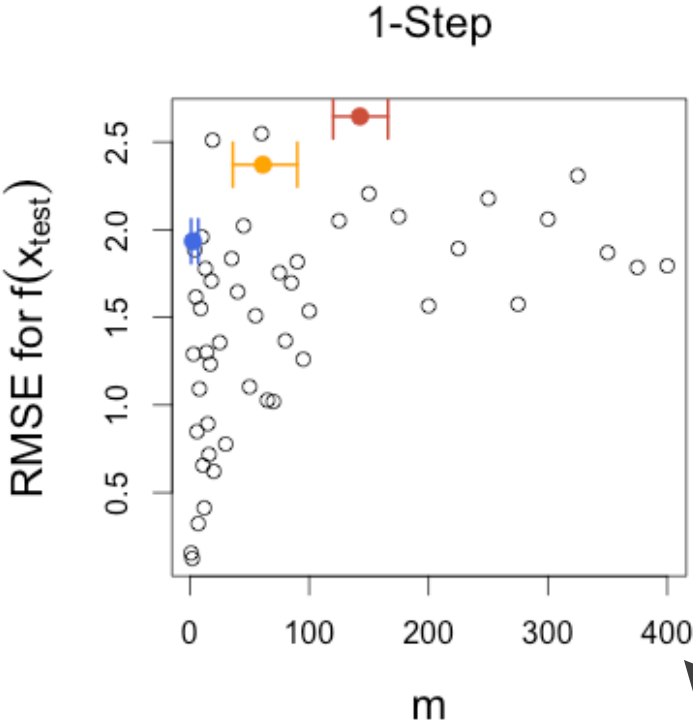
300-Step

Welch



- $\kappa_0 = 3$
- $\kappa_0 = 100$
- $\kappa_0 = \infty$
- Fixed  $m$

# Results



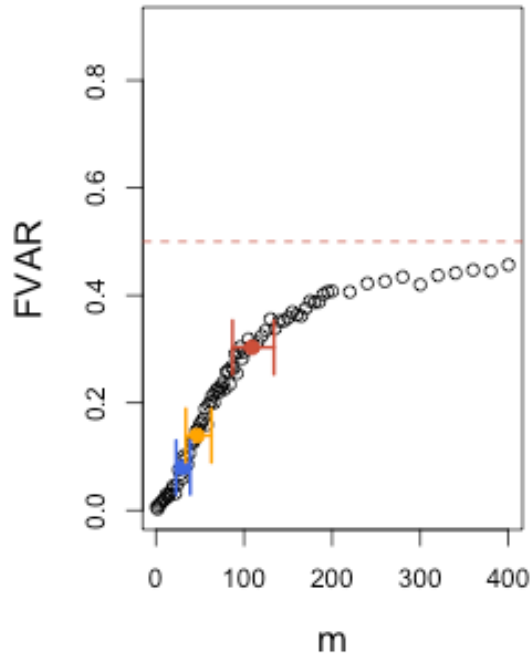
$p = 15$  inputs (each a different branch)  
 $f$  not additive at all!

- $\kappa_0 = 3$
- $\kappa_0 = 100$
- $\kappa_0 = \infty$
- Fixed  $m$

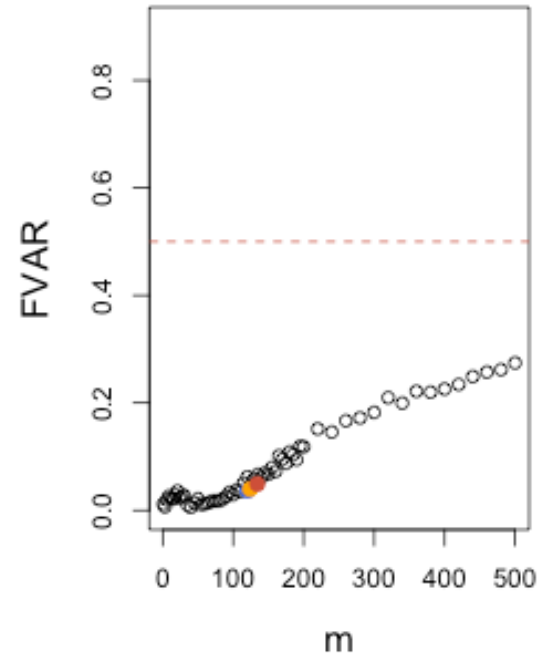
1-step:  $f(x) = 20I(x \geq 0.5)$  ( $p = 1$ )

# Variable Selection

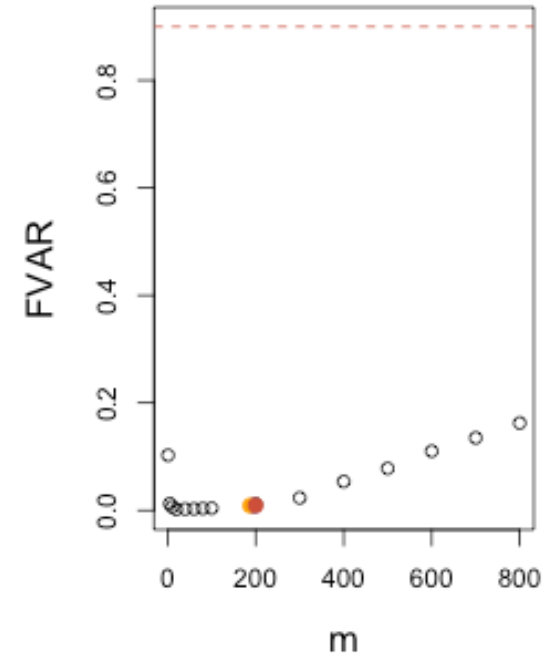
Friedman:  $p = 10$  (5 real)



100-Step:  $p = 200$  (100 real)



Welch:  $p = 200$  (20 real)



- $\kappa_0 = 3$
- $\kappa_0 = 100$
- $\kappa_0 = \infty$
- Fixed  $m$

FVAR = Proportion of branches involving “false” variables



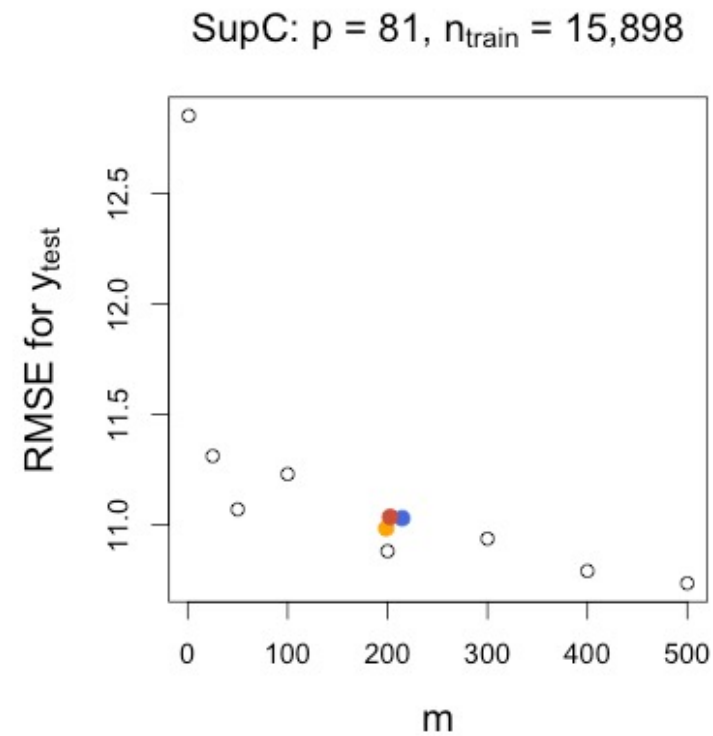
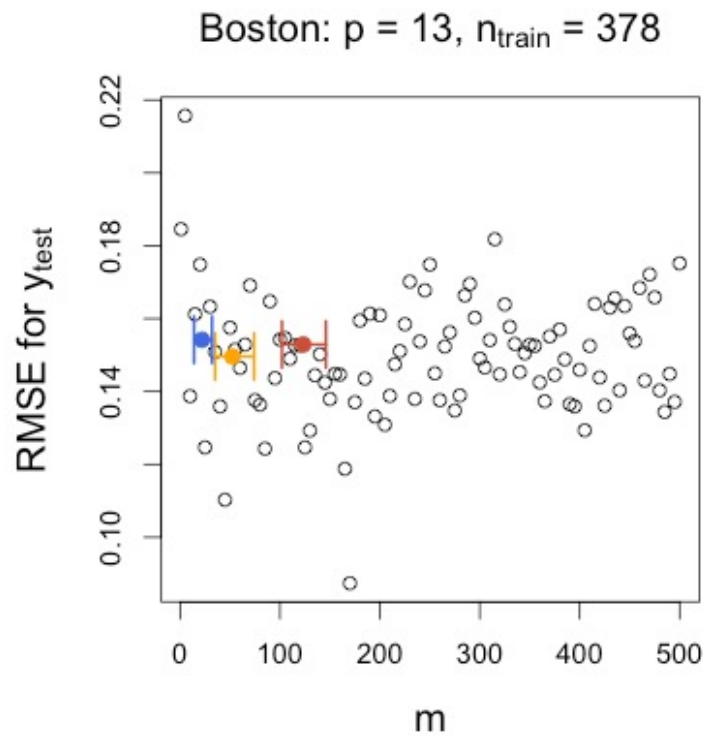
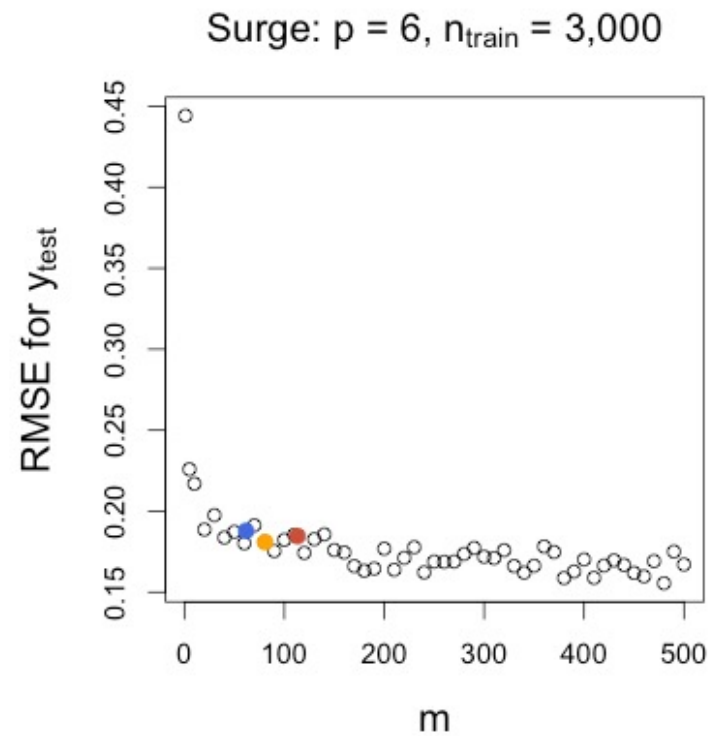
1. Recap of BART
- 2. Bayesian Inference of the Number of Trees**
  - i. Motivation
  - ii. A Fully Bayesian Model
  - iii. Sampling from the Posterior Distribution
  - iv. Code
  - v. Simulations
  - vi. Application to Real Data**
3. Conclusion

# Real Datasets

## Real Data

Dataset	$n_{\text{train}}$	$n_{\text{test}}$	$p$	$N_{\text{mc}}$ (infer $m$ )	$N_{\text{mc}}$ (fix $m$ )
Surge	3,000	1,000	6	62,000 (✗)	22,000
Boston	378	128	13	1,000,000 (✓)	3,000
Superconductor	15,898	5,299	81	100,000 (✗)	100,000

# Results



- $\kappa_0 = 3$
- $\kappa_0 = 100$
- $\kappa_0 = \infty$
- Fixed  $m$

1. Recap of BART
2. Bayesian Inference of the Number of Trees
  - i. Motivation
  - ii. A Fully Bayesian Model
  - iii. Sampling from the Posterior Distribution
  - iv. Code
  - v. Simulations
  - vi. Application to Real Data
3. Conclusion

# Conclusions

- Bayesian Inference of  $m$  generally works well
  - Accurate predictions
  - Variable selection
  - Convenience

# Conclusions

- Bayesian Inference of  $m$  generally works well
  - Accurate predictions
  - Variable selection
  - Convenience
  - Sometimes underfit ( $m$  too small)
  - Computation time?

# Conclusions

- Bayesian Inference of  $m$  generally works well
  - Accurate predictions
  - Variable selection
  - Convenience
  - Sometimes underfit ( $m$  too small)
  - Computation time?
- **Just fix  $m = 200$ ?**

# Conclusions

- Bayesian Inference of  $m$  generally works well
  - Accurate predictions
  - Variable selection
  - Convenience
  - Sometimes underfit ( $m$  too small)
  - Computation time?
- Just fix  $m = 200$ ?
- Recommendation: Truncated Poisson prior with  $\theta = 200$  ( $\kappa_0 = \infty$ )



# Conclusions

- Bayesian Inference of  $m$  generally works well
  - Accurate predictions
  - Variable selection
  - Convenience
  - Sometimes underfit ( $m$  too small)
  - Computation time?
- Just fix  $m = 200$ ?
- Recommendation: Truncated Poisson prior with  $\theta = 200$  ( $\kappa_0 = \infty$ )
- Maybe also try  $\kappa_0 = 3$

# Conclusions

- Bayesian Inference of  $m$  generally works well
  - Accurate predictions
  - Variable selection
  - Convenience
  - Sometimes underfit ( $m$  too small)
  - Computation time?
- Just fix  $m = 200$ ?
- Recommendation: Truncated Poisson prior with  $\theta = 200$  ( $\kappa_0 = \infty$ )
- Or try two values of  $\kappa_0$
- **Boosting**

Thank You!