

# **Bayesian Additive Regression Trees**

**Lecture 23 - CS 689, Spring 2023**

- Last time: decision trees and their ensembles
  - 1. Individual decision trees
  - 2. DTree ensemble with random forests
  - 3. DTree ensembles with greedy optimization - forward stagewise learning and **gradient boosted trees** (next slide)
- Rest of today: **Bayesian learning** for decision tree ensemble

$$L = (y - f(x))^2$$

$$\frac{\partial L}{\partial f} = -2(y - f(x))$$

---

### Algorithm 16.4: Gradient boosting

---

- 1 Initialize  $f_0(\mathbf{x}) = \operatorname{argmin}_{\gamma} \sum_{i=1}^N L(y_i, \phi(\mathbf{x}_i; \gamma))$ ;
  - 2 **for**  $m = 1 : M$  **do**
  - 3     Compute the gradient residual using  $r_{im} = - \left[ \frac{\partial L(y_i, f(\mathbf{x}_i))}{\partial f(\mathbf{x}_i)} \right]_{f(\mathbf{x}_i) = f_{m-1}(\mathbf{x}_i)}$
  - 4     Use the weak learner to compute  $\gamma_m$  which minimizes  $\sum_{i=1}^N (r_{im} - \phi(\mathbf{x}_i; \gamma_m))^2$ ;
  - 5     Update  $f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \nu \phi(\mathbf{x}; \gamma_m)$ ;
  - 6 Return  $f(\mathbf{x}) = f_M(\mathbf{x})$
- 

learning rate / step size

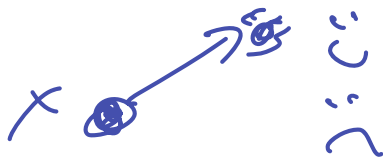
# Review: MCMC for Bayesian pred.

- Goal: want posterior predictive

$$p(y/x, x^{\text{Tr}}, y^{\text{Tr}}) = \int_{\Theta} \underbrace{p(\theta/x^{\text{Tr}}, y^{\text{Tr}})}_{\text{posterior}} \underbrace{p(y/x, \theta)}_{\text{likelihood}} d\theta$$
$$\approx \sum_{s=1}^S p(y/x, \theta^{(s)}) \quad \text{for } \theta^{(s)} \sim p(\theta/x^{\text{Tr}}, y^{\text{Tr}})$$

- Method: Approximate with samples from the parameter posterior
  - No closed form, but can use Markov Chain Monte Carlo

# MCMC: key algorithms



## Algorithm 24.2: Metropolis Hastings algorithm

- 1 Initialize  $x^0$  ;
- 2 **for**  $s = 0, 1, 2, \dots$  **do**
- 3     Define  $x = x^s$  ;
- 4     Sample  $x' \sim q(x'|x)$  ;
- 5     Compute acceptance probability

$$\alpha = \frac{\tilde{p}(x')q(x|x')}{\tilde{p}(x)q(x'|x)}$$

- 6     Compute  $r = \min(1, \alpha)$  ;
- 7     Sample  $u \sim U(0, 1)$  ;
- 7     Set new sample to

$$x^{s+1} = \begin{cases} x' & \text{if } u < r \\ x^s & \text{if } u \geq r \end{cases}$$

## Algorithm 1: Gibbs sampler

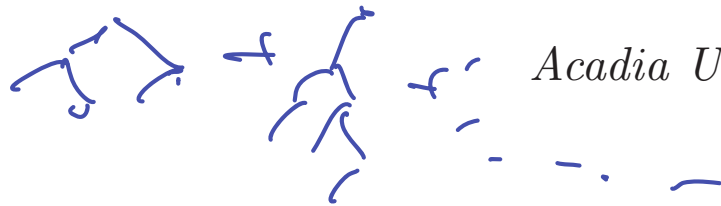
- Initialize  $x^{(0)} \sim q(x)$
- for** iteration  $i = 1, 2, \dots$  **do**
- $x_1^{(i)} \sim p(X_1 = x_1 | X_2 = x_2^{(i-1)}, X_3 = x_3^{(i-1)}, \dots, X_D = x_D^{(i-1)})$
  - $x_2^{(i)} \sim p(X_2 = x_2 | X_1 = x_1^{(i)}, X_3 = x_3^{(i-1)}, \dots, X_D = x_D^{(i-1)})$
  - $\vdots$
  - $x_D^{(i)} \sim p(X_D = x_D | X_1 = x_1^{(i)}, X_2 = x_2^{(i)}, \dots, X_{D-1} = x_{D-1}^{(i)})$
- end for**

for us

$$\tilde{p}(\theta) = \underbrace{p(\theta)}_{\text{prior}} \underbrace{p(Y/x, \theta)}_{\text{lik}}$$

# **BART: BAYESIAN ADDITIVE REGRESSION TREES<sup>1,2</sup>**

BY HUGH A. CHIPMAN, EDWARD I. GEORGE AND ROBERT E. MCCULLOCH



*Acadia University, University of Pennsylvania and  
 University of Texas at Austin*

- **Model:**

$$\text{ensemble } f(x) = \sum_{j=1}^m \underbrace{g_j(x)}_{\rightarrow \text{one tree}}$$

$P(g)$ : prior on trees

$P(y/x, f)$ : likelihood,  $N(f(x), \sigma^2)$

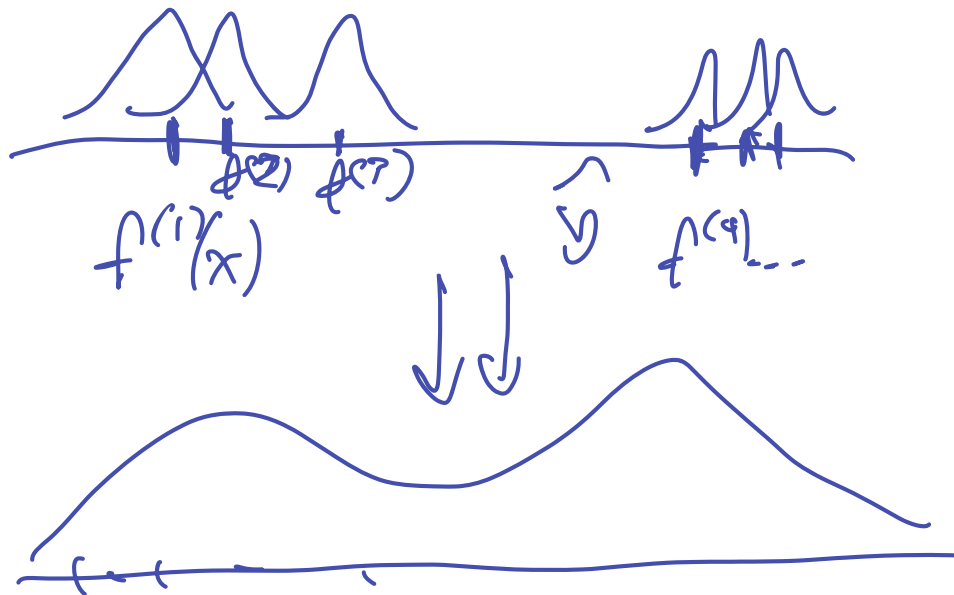
Posterior over ensembles:  $P(f/y^{tr}, y^{tr}) \propto P(f) P(y^{tr}/x^{tr}, f)$

# Uncertainty-aware predictions

BART for causal inference

- If we had samples of posterior ensembles, we could do lots of nice probabilistic inferences / predictions!

$$f^{(s)} \sim p(f/x, \tau)$$



$$P(Y \leq t/x) \approx$$

$$\frac{1}{S} \sum_s P(Y \leq t / f^{(s)}(x))$$

# Model and Prior Structure

$T$ : binary tree,  $M = \{M_1, \dots, M_b\}$  for each leaf

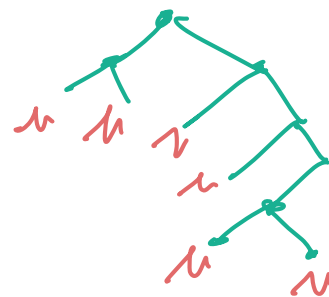
File  $Y = \sum_{j=1}^m g(x; \underset{\uparrow}{T_j}, \underset{\uparrow}{M_j}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$

$$P((T_1, M_1) \dots (T_m, M_m), \sigma) = P(\sigma) \left[ \underbrace{T_j}_{\text{leaf preds.}} \underbrace{P(M_j | T_j)}_{\text{struct}} \right]$$

Each tree is indep

$$p(\mathbf{M}_j | \mathbf{T}_j) = \prod_{i=1}^b p(M_{ij} | T_j)$$

each col indep





# Tree structure prior

Used to regularize size/shape of trees

• Node at depth  $d$ , is non-terminal with prob

$$\propto (1+d)^{-\beta}$$

↑                      ↑  
hyperparams

- Uniform dist. b.s.:
- ① What variable to split on
  - ② What threshold to use:  
any discrete value in training set

# Data-driven priors

$$(u_{ij} | T_j) \sim N(\underline{\mu}_u, \sigma_u^2)$$

- Leaf predictions: conjugate normal prior

1-D Gaussian priors

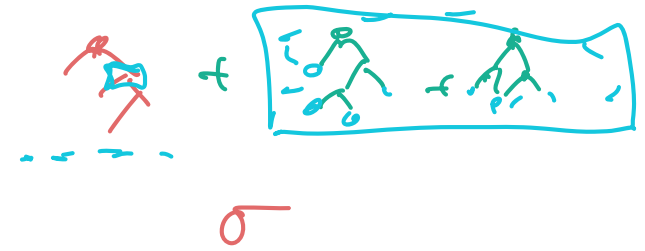
$$p(y/x \in \text{leaf } b) = \int_{\mathcal{M}_b} p(\mathcal{M}_b/T) p(y^{\text{Tr}}/x^{\text{Tr}}, \mathcal{M}_b) \\ = N(\text{---}, \text{---})$$

- Variance term: conjugate inverse chi square

Avg of  $\mathcal{M}_b$ ,  
and avg  $y$  within leaf  $b$

$\mu_u, \sigma_u^2$  = based on  $y^{\text{Tr}}$

# Backfitting MCMC



- Gibbs sampler to resample one tree at a time, and variance

$$P(\underbrace{(T_1, M_1), \dots, (T_m, M_m)}_{\text{trees}}, \sigma | X, Y)$$

GS for  $P(\sigma | T_1, M_1, \dots, T_m, M_m, X, Y)$

then each  $P(\underbrace{T_j, M_j}_{\text{tree } j} | \underbrace{T_{(j)}, M_{(j)}}_{\text{other trees}}, X, Y, \sigma)$

# Partial residuals



$$\vec{R}_j \equiv \left\{ y_i - \sum_{k \neq j} g(x_i; T_k, M_k) \right\}_{i=1 \dots N}$$

$$p(T_j, M_j \mid T_{(j)}, M_{(j)}, X, Y, \sigma) = p(T_j, M_j \mid R_j, \sigma)$$

$$\propto \underbrace{p(R_j \mid T_j, M_j, \sigma) p(T_j, M_j)}$$

$$R_j = g(x; T_j, M_j) + \varepsilon$$

# Tree proposal & resampling

$$\Rightarrow p(\tau_j | R_j, \sigma) \propto p(\tau_j) \underbrace{\int_{M_j} p(R_j | M_j, \tau_j, \sigma) p(M_j | \tau_j, \sigma) dM_j}_{\text{Gaussian } R_j \sim \mathcal{N}(\dots, \dots)}$$

- Proposal distribution from old to new tree:

$$\Rightarrow q(\tau^{\text{new}} | \tau^{\text{old}})$$

→ Grow terminal node



→ Prune pair of nodes



→ Change nonAvar. rule



$$\alpha = \frac{p(\tau^{\text{new}} | R_j, \sigma) q(\tau^{\text{old}} | \tau^{\text{new}})}{p(\tau^{\text{old}} | R_j, \sigma) q(\tau^{\text{new}} | \tau^{\text{old}})}$$

# Leaf resampling

For each iteration,  $S \sim$

for each tree  $j$ ,

$$(T_j^{(s)}, M_j^{(s)}) \sim \text{MCMC acceptance for } P(T, M; R_j, \sigma)$$

$$\sigma^{(s)} \sim P(\sigma | T, M, \sigma, X, Y) \text{ [conjugate]}$$

$\Rightarrow$

$$\left\{ T_1^{(s)}, M_1^{(s)}, \dots, T_m^{(s)}, M_m^{(s)}, \sigma^{(s)} \right\}_{s=1 \dots S}$$









