Multidimensional Monotonicity Discovery with mBART

Robert McCulloch

ASU Rober.McCulloch@asu.edu

Collaborations with: H. Chipman, **E. George**, T. Shively

Plan

- ► I) BART (Bayesian Additive Regression Trees)
- ► II) Monotone BART: mBART
- III) Monotonicity Discovery with mBART

Part I. BART (Bayesian Additive Regression Trees)

The Fundamental Regression Setup:

• Data: *n* observations of *y* and $x = (x_1, ..., x_p)$

Suppose:
$$Y = f(x) + \epsilon$$
, $\epsilon \sim N(0, \sigma^2)$

Bayesian Ensemble Idea: Approximate unknown f(x) by the form

$$f(x) = g(x; \theta_1) + g(x; \theta_2) + \dots + g(x; \theta_m)$$
$$\theta_1, \theta_2, \dots, \theta_m \quad \text{iid} \sim \pi(\theta)$$

and use the posterior of f given y and x for inference.

BART: Each $g(x; \theta_j)$ is a regression tree function. Key calibration: Using y, set $\pi(\theta)$ so that $Var(f) \approx Var(y)$.

Beginning with a Single Tree Model

Let *T* denote the tree structure including the decision rules

Let $M = {\mu_1, \mu_2, \dots, \mu_b}$ denote the set of bottom node μ 's.

Let g(x; T, M)be a regression tree function that assigns a μ value to x



A single tree model:

 $Y = g(x; T, M) + \sigma z, z \sim N(0,1)$

 $"\theta = (T, M)"$

Bayesian CART: Just add a prior $\pi(M, T)$

Bayesian CART Model Search Chipman, George, McCulloch (1998 JASA), Mallick and Smith (1998)

 $\pi(M, T) = \pi(M \mid T)\pi(T)$

 $\pi(T)$: Stochastic process to generate tree skeleton plus uniform prior on splitting variables and splitting rules.

$$\pi(M \mid T) : (\mu_1, \mu_2, \ldots, \mu_b)' \sim N_b(0, \tau^2 I)$$

Closed form for $\pi(T | y)$ facilitates MCMC stochastic search for promising trees.

Moving on to BART

Bayesian Additive Regression Trees Chipman, George, McCulloch (2007 NIPS, 2010 AOAS)

The BART ensemble model

 $Y = g(x; T_1, M_1) + g(x; T_2, M_2) + \ldots + g(x; T_m, M_m) + \sigma z, \quad z \sim N(0, 1)$



Each (T_i, M_i) identifies a single tree.

For each x, Y is the sum of m bottom node μ 's, plus noise.

Number of trees m can be much larger than sample size n.

 $g(x; T_1, M_1), g(x; T_2, M_2), ..., g(x; T_m, M_m)$ is a highly redundant "over-complete basis" with many many parameters.

Complete the Model with a "Regularization/Boosting" Prior

$$\pi((T_1, M_1), (T_2, M_2), \dots, (T_m, M_m), \sigma)$$

 π applies the Bayesian CART prior to each (T_j, M_j)

- Hyperparameters of π are set so that both the T_j 's and the μ_{ij} 's are likely to be small.
- Thus, π keeps the contribution of each $g(x; T_j, M_j)$ small, to explain only a small portion of the fit.

The observed variation of Y is used to guide the hyperparameter settings for the μ and σ priors.

• Note that because
$$\mu_{ij}$$
's iid $\sim N(0, \tau^2)$, we have

$$Var(Y) \approx Var(\Sigma_{ij}\mu_{ij}) = m\tau^2,$$

which facilitates the calibration of au^2 for BART's success!

7 / 50

Build up the fit, by adding up tiny bits of fit ...



Build up the fit, by adding up tiny bits of fit ...



Connections to Other Modeling Ideas

$$\begin{split} \mathsf{Y} &= \mathsf{g}(\mathsf{x};\mathsf{T}_1,\mathsf{M}_1) + \ldots + \mathsf{g}(\mathsf{x};\mathsf{T}_\mathsf{m},\mathsf{M}_\mathsf{m}) + \sigma \, \mathsf{z} \\ & \mathsf{plus} \\ \pi((\mathsf{T}_1,\mathsf{M}_1),\ldots,(\mathsf{T}_\mathsf{m},\mathsf{M}_\mathsf{m}),\sigma) \end{split}$$

Bayesian Nonparametrics:

- Lots of parameters (to make model flexible)
- A strong prior to shrink towards simple structure (regularization)
- BART shrinks towards additive models with some interaction

Gradient Boosting:

Fit becomes the cumulative effort of many weak learners

Dynamic Random Basis Elements:

• $g(x; T_1, M_1), ..., g(x; T_m, M_m)$ are dimensionally adaptive

Some Distinguishing Features of BART

$$\begin{split} \mathsf{Y} &= \mathsf{g}(\mathsf{x};\mathsf{T}_1,\mathsf{M}_1) + \ldots + \mathsf{g}(\mathsf{x};\mathsf{T}_m,\mathsf{M}_m) + \sigma \, \mathsf{z} \\ & \mathsf{plus} \\ \pi((\mathsf{T}_1,\mathsf{M}_1),\ldots,(\mathsf{T}_m,\mathsf{M}_m),\sigma) \end{split}$$

BART is NOT Bayesian model averaging of a single tree model

Unlike boosting and random forests, the BART algorithm updates a fixed set of m trees, over and over

Choose *m* large for best estimation of E[Y|x] and prediction

More trees yields more approximation flexibility

Choose m small to measure variable importance

Fewer trees force the x's to compete for entry

A Sketch of the BART MCMC Algorithm

$$\begin{split} Y &= g(x; T_1, M_1) + ... + g(x; T_m, M_m) + \sigma \, z \\ & \text{plus} \\ \pi((T_1, M_1), ..., (T_m, M_m), \sigma) \end{split}$$

Bayesian Backfitting: Outer loop is a "simple" Gibbs sampler

$$(T_i, M_i) \mid Y$$
, all other (T_j, M_j) , and σ
 $\sigma \mid Y, (T_1, M_1, \dots, T_m, M_m)$

To draw (T_i, M_i) above, subtract the contributions of the other trees from both sides to get a simple one-tree model.

We integrate out M_i to draw T_i and then draw $M_i | T_i$.

... as the MCMC runs, trees in the sum will grow and shrink, swapping fit amongst them

Using the MCMC Output to Draw Inference

Each iteration d yields a new draw from the posterior of f

$$\hat{f}_d(\cdot) = g(\cdot; T_{1d}, M_{1d}) + \cdots + g(\cdot; T_{md}, M_{md})$$

To estimate f(x) we simply average the $\hat{f}_d(\cdot)$ draws at x

Posterior uncertainty is captured by the variation of the $\hat{f}_d(x)$ eg, 95% credible region estimated by middle 95% of values

Can do the same with functionals of f.

Automatic Uncertainty Quantification

A simple simulated 1-dimensional example

95% pointwise posterior intervals, BART



Out of Sample Prediction Comparisons

Predictive comparisons on 42 data sets

Data from Kim, Loh, Shih and Chaudhuri (2006) (thanks Wei-Yin Loh!)

- ▶ p = 3 to 65, n = 100 to 7,000.
- ▶ For each data set 20 random splits into 5/6 train and 1/6 test.
- Use 5-fold CV on train to pick hyperparameters (except BART-default!).
- Gives 20*42 = 840 out-of-sample predictions. For performance comparisons, the RMSE of each method is divided by the smallest RMSE of all.

- + Each boxplot represents 840 predictions for a method
- + 1.2 means you are 20% worse than the best
- + BART-cv best
- + BART-default does amazingly well!!



Measuring Variable Importance with BART

When $Y = g(x;T_1,M_1) + ... + g(x;T_m,M_m) + \sigma z$ is fit to data, we can count how many times a predictor is used in the trees.

For example, in the tree here, x_2 and x_5 are each used once.

The importance of each \underline{x}_k can thus be measured by its overall usage frequency.

This approach is most effective when the number of trees <u>m</u> is small.



A Competitive Bottleneck for Entry



BART is an automatic attractor of x's to explain Y

Those x's which most explain Y are attracted most

Small number of trees \underline{m} creates a bottleneck which excludes useless $\underline{x's}$

Example: The Friedman Test Function

BART variable importance on data simulated from:

 $Y = 10\sin(\pi x_1 x_2) + 20(x_3 - .5)2 + 10x_4 + 5x_5 + 0x_6 + \dots + 0x_{10} + \epsilon$

Variable usage frequencies as the number of trees m is reduced



mBART: Multidimensional Monotone BART Chipman, George, McCulloch, Shively (2021 Bayesian Analysis)

The Key Idea:

BART approximates a function by a sum of tree functions

mBART approximates a monotone function by a sum of monotone tree functions

This works because of the obvious fact:

The sum of monotone functions yields a monotone function

An Example of a Monotone Tree Function







Three different views of a bivariate monotone tree.

In what sense is this tree function monotone?



A tree function g(x; T, M) is said to be monotone nondecreasing in x_i if for all x_{-i} and $\delta > 0$,

$$g(x_i, x_{-i}; T, M) \leq g(x_i + \delta, x_{-i}; T, M)$$

For simplicity and wlog, let's restrict attention to monotone nondecreasing tree functions.

The mBART Prior

```
\begin{split} Y &= g(x; T_1, M_1) + ... + g(x; T_m, M_m) + \sigma \, z \\ & \text{plus} \\ \pi((T_1, M_1), ..., (T_m, M_m), \sigma) \end{split}
```

Recall the BART parameter

$$\theta = ((T_1, M_1), (T_2, M_2), \dots, (T_m, M_m), \sigma)$$

Let $S = \{\theta : each (T_j, M_j) \text{ is monotone in a desired subset of } x'_i s \}$

To impose the monotonicity we simply truncate the BART prior $\pi(\theta)$ to the set S

 $\pi^*(heta) \propto \pi(heta) I_{\mathcal{S}}(heta)$

where $I_{S}(\theta)$ is 1 if *every* tree in θ is monotone.

Forcing a tree to be monotone is easy: we simply constrain the mean level of a node to be greater than those of its "below-neighbors", and less than those of its "above-neighbors".



For example, the mean level of node 13 must be greater than those of 10 and 12 and less than that of node 7.

For any bottom node μ , given the rest of the tree, we can figure out (and easily code) its interval of constraint.

Because we only make local changes via the MCMC algorithm, this criterion suffices for all computations.

The remaining challenge is the construction of a new algorithm which can handle the nonconjugacy of truncated priors on μ 's.

A New BART MCMC "Christmas Tree" Algorithm

 $\pi((T_1, M_1), (T_2, M_2), \dots, (T_m, M_m), \sigma | y))$

Bayesian Backfitting again: Iteratively sample each (T_j, M_j) given (y, σ) and other (T_j, M_j) 's

Each $(T^0, M^0) \rightarrow (T^1, M^1)$ update is sampled as follows:

Only $M^0_{Old} \rightarrow M^1_{New}$ needs to be updated.

Works for both BART and mBART.

Example: Product of two x's

Let's consider a very simple simulated monotone example:

 $Y = x_1 x_2 + \epsilon$, $x_i \sim \text{Uniform}(0, 1)$.

Here is the plot of the true function $f(x_1, x_2) = x_1 x_2$



First we try a single (just one tree), unconstrained tree model.

Here is the graph of the fit.



The fit is not terrible, but there are some aspects of the fit which violate monotonicity.

Here is the graph of the fit with the monotone constraint:



We see that our fit is monotonic, and more representative of the true f.

Here is the unconstrained BART fit:



Much better (of course) but not monotone!

And, finally, the constrained BART fit:



Not Bad!

Same method works with any number of x's!

Automatic Uncertainty Quantification

Revisiting our simple simulated 1-dimensional example



mBART intervals are tighter!

Example: RMSE Reduction by Monotone Regularization

$$Y = x_1 x_2^2 + x_3 x_4^3 + x_5 + \epsilon,$$

$$\epsilon \sim N(0, \sigma^2), \ x_i \sim \text{Uniform}(0, 1).$$

For various values of σ , we simulated 5,000 observations.

RMSE improvement of mBART over BART

σ	Monotone BART	Unconstrained BART	Percentage
	RMSE	RMSE	Increase
0.5	0.14	0.16	14%
1.0	0.17	0.28	65%



 $\sigma = 0.2, 0.5, 0.7, 1.0$

Suppose we don't know if f(x) is monotone up, monotone down or even monotone at all.

Of course, a simple strategy would be to simply compare the fits from BART and mBART.

Good news! We can do even better than this by deploying mBART to simultaneously estimate all the monotone components of f.

With this strategy, monotonicity can be discovered rather than imposed!

To begin simply, suppose x is one-dimensional and f is of bounded variation.

The Jordan Decomposition Theorem: Any such f can be uniquely written (up to an additive constant) as the sum of its monotone up and monotone down components

$$f(x) = f_{up}(x) + f_{down}(x)$$

where

- when f(x) is increasing, f_{up}(x) increases at the same rate and is flat otherwise,
- when f(x) is decreasing, f_{down}(x) decreases at the same rate and is flat otherwise.

The Monotone Discovery Strategy with mBART

Key Idea: To discover the monotone decomposition of f, we treat f(x) as embedded in a two-dimensional function

$$f^*(x_1, x_2) = f_{up}(x_1) + f_{down}(x_2).$$

Letting $x_1 = x_2 = x$ be duplicate copies of x, we simply estimate $f^*(x_1, x_2)$ with mBART

- constrained to be monotone up in the x_1 direction, and
- constrained to be monotone down in the x₂ direction.

Thus, we are estimating the monotone "projections" of $f^*(x_1, x_2)$ along the x_1 and x_2 axes, i.e.

•
$$P_{[x_1]}f^*(x_1, x_2) = f_{up}(x_1)$$

• $P_{[x_2]}f^*(x_1, x_2) = f_{down}(x_2)$

Example: Suppose $Y = x^3 + \epsilon$.



BART and mBART

mBARTD, fup, fdown

Note that $\hat{f}_{down} \approx 0$ (the red in the right plot), as we would expect when f is monotone up.

Remark: mBARTD = $\hat{f}_{up} + \hat{f}_{down}$ is an alternative estimate of f

As the sample size is increased from 200 to 1,000, \hat{f}_{down} gets even flatter.



Suggests consistent estimation of the monotone components!!

Example: Suppose $Y = x^2 + \epsilon$.



BART and mBART

mBARTD, fup, fdown

On the left, BART is good, but simple mBART is not.

- On the right, \hat{f}_{up} and \hat{f}_{down} are spot on.
- And mBARTD = $\hat{f}_{up} + \hat{f}_{down}$ seems even better than BART!

Example: Suppose
$$Y = sin(x) + \epsilon$$
.



- BART is great, but simple mBART reveals nothing.
- \hat{f}_{up} and \hat{f}_{down} have discovered the monotone decomposition.
- And mBARTD = $\hat{f}_{up} + \hat{f}_{down}$ is great too.

To extend this approach to multidimensional x, we simply duplicate each and every component of x !!!

Example: House Price Data

n = 128 houses, y = house price (\$ thousands), x = (nbhd (1,2 or 3), size (sq ft thousands), brick (B or N)).Call lm(formula = price ~ nbhd + size + brick, data = hdat) Coefficients: Estimate Std. Error t value Pr(>|t|) (Intercept) 18.725 10.766 1.739 0.0845 . nbhd2 5 556 2 779 1 999 0 0478 * 36.770 2.958 12.430 < 2e-16 *** nbhd3 size 46.109 5.527 8.342 1.25e-13 *** 2.438 7.855 1.69e-12 *** brickYes 19.152 _ _ _

Residual standard error: 12.5 on 123 degrees of freedom Multiple R-squared: 0.7903,Adjusted R-squared: 0.7834 F-statistic: 115.9 on 4 and 123 DF, p-value: < 2.2e-16

If the linear model is correct, we are monotone up in all three variables.

Remark: For the linear model we have to dummy up *nbhd*, but for BART and mBART we can simply leave it as an ordered numerical categorical variable.

Let's first compare BART, mBART (constrained up), and mBARTD to estimate the effect of *size* conditionally on the six possible values of (*nbdh*, *brick*)



Note how mBARTD = $\hat{f}_{up} + \hat{f}_{down}$ adaptively shrinks the estimates towards the mBART estimates.

The full picture emerges from estimates of the effect of *size* via \hat{f}_{up} and \hat{f}_{down} conditionally on the six possible values of (*nbdh*, *brick*)



Price is clearly conditionally monotone up in all three variables!

By simultaneously estimating \hat{f}_{up} and \hat{f}_{down} , we have discovered monotonicity without any imposed assumptions!!!

This can all be most conveniently done using the mBART variable importance strategy to gauge the relationships between y = price and x = (sizeUp, sizeDn, nbhdUp, nbhdDn, brickUp, brickDn)y = price and x = (UpPrick, UpNhbd, UpSize, DpPrick, DpNhbd, DpSize)

x = (UpBrick, UpNbhd, UpSize, DnBrick, DnNbhd, DnSize).

x = (sizeUp, sizeDn, nbhdUp, nbhdDn, brickUp, brickDn).

This frequency-of-use variable importance strategy reveals clearly that *price* is conditionally monotone up in all three variables:





Posterior distribution of percent of rules in tree ensemble using a variable

Example: The Diabetes Data

Benchmark Dataset used in *Least Angle Regression* Efron, Hastie, Johnstone, Tibshirani (2004, *AOS*)

n = 442 diabetes patients, y = disease progression measure, x = (age, bmi, glu, hdl, ldl, ltg, map, sex, tc, tch)

The mBART variable importance strategy identifies six important variables together with the direction of their conditional effects



bmi-U, Itg-U, map-U, glu-U, hdl-D, sex-D.

- Discovery variable importance: put in x and -x.
- BART Prior: put tight prior on σ saying you want the same kind of σ as you got from classic BART.



bmi-U, Itg-U, map-U, glu-U, hdl-D, sex-D.

Coefficients: Estimate Std. Error t value Pr(>|t|) (Intercept) -3.710e-08 2.576e+00 0.000 1.000000 -4.768e-01 2.845e+00 -0.168 0.867000 age -1.142e+01 2.915e+00 -3.917 0.000104 *** sex 2.475e+01 3.168e+00 7.813 4.30e-14 *** bmi map 1.545e+01 3.115e+00 4.958 1.02e-06 *** tc -3.772e+01 1.984e+01 -1.901 0.057947 . 1d1 2.270e+01 1.614e+01 1.406 0.160389 hd1 4.812e+00 1.012e+01 0.475 0.634720 tch 8.432e+00 7.689e+00 1.097 0.273456 3.578e+01 8.186e+00 4.370 1.56e-05 *** ltg 3.220e+00 3.142e+00 1.025 0.305998 glu ---Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 54.15 on 431 degrees of freedom Multiple R-squared: 0.5177,Adjusted R-squared: 0.5066 F-statistic: 46.27 on 10 and 431 DF, p-value: < 2.2e-16

Not the same !!!!

which one is more useful !!!????

How about this ????

	Estimate	Int. Secur	a value	Prorinti
(Intercept)	-3.730+18	2.532+00	4.688	1.0000
@	2.411+10	3.120+00	4.774	0.6883
	-1.273a-01	3.1084-00	-6.006	6. ifa-if. ***
hai.	2.064400	4.000+00	R. 696	0.324-68 ***
***	1.432+11	3.600400	6.736	3. ilu-id +++
**	-6.767#902	2.885.4413	-0.060	0.9624
145	1.4444-02	2.6364-63	0.007	0.9524
hall.	£.00+01	1.076-03	0.049	0.0011
***	3.0684-00	1.3634411	0.20	0.7860
248	8.784+01	0.000-02	0.00	0.4268
gin .	3.088.00	3.900-00	0.884	0.2733
*gx. *	1.000000	1.1842400		0.000
	a starter	3.000200		0.0000
10.2	3.1714-12	3.364-12	0.945	0.005
16.2	1.706+12	2.536+12	4.03	0.5016
hil. 2	8.000-00	7.0764-04	1.000	0.2770
wah. 2	3.6834-01	2.800+11	1.2%	0.2034
14g.2	6.034-03	8.200-01	4.838	0.000
gin 2	1.434+10	4.400-00	1.20	0.2260
4ge. 161	7.000+00	3.4864-00	2.635	0.0435.*
age, bui	-0.004-01	3.7864-00	-0.227	0.8208
40.00	8.826-01	3.435+00	0.243	0.882
age. to	-7.5864-00	2.000+01	-4.90	0.7969
age 141	-3.204+00	2.355a-01	-9.106	0.0000
age hill	8.054-00	1.3364-01	4.746	0.4963
age, tok	8.856+10	1.000+01	4.879	0.2798
age ing	5.657+100	1.0854-01	4.607	0.6778
aga gin	2.000+00	3.827+100	4.778	0.6967
san bei	3.077+100	3.7104-00	0.828	0.4074
san nap	6.063a-00	3.530+00	1.184	0.2273
ware. No	2.084+01	2.8534-04	0.794	0.004
aan. 141	-1.604-01	2.250+14	-9.782	0.003
earn halfs	-5.9604-00	1.3564-01	-0.488	0.6211
ware, took	-6.2884-00	0.5304-00	-9.607	0.6005
san. ing	-5.486+100	1.0794-01	-0.626	0.5896
san gin	2.1784-00	3.6274-00	0.631	0.6348
hai.mp	7.3684-00	4.1114-00	1.792	0.6739 .
loni. to	-1.684-01	3.1814-01	-0.00	0.4514
bai. 141	1.100-01	2.672/414	0.65	0.6070
hai. Mi	L. ACT #-00	1.014-01	0.270	0.7018
hai. tok	-1.000400	1.0004-01	-0.046	0.0000
hai. ing	1.4114-00	1.200401	0.000	0.0044
ini, gin	1.1134-00	4.356,000	0.20	0.7676
map. to	3.284-01	3.269-01	0.764	0.007
sap. one		2.100251	-0.668	0.0000
sup su	-	1.000000	-0.646	0.0000
sap ton	-2.110.000	1.000000		0.0004
	-	4.348-000		0.007
	-1.630-000	1.414-012		0.000
		1.017-017		0.000
to.tok	-1.000+02	8.300+11	-1.20	0.2013
weing	-1.810+12	4.2754-62	-1.200	0.7730
to gin	-1.384-10	2.836-14	-0.206	4.7673
14.66	1.258+12	1.00440	4.655	0.4945
141.144	6.7674-01	7.000+04	0.831	0.604
141. ing	1.320+12	6.00440	4.203	0.0014
14. gin	4.072+00	2.435a-01	4.05	0.8655
hill, tak	5.684-01	4.7754-94	1.186	0.2365
hill-ing	6.0884-01	2.18fa=02	0.348	0.7884
hill gin	1.6Mar01	1.634-01	0.733	0.0000
wah. ing	1.854+11	2.075a-01	0.634	0.6390
ush gin	1.1224-01	1.120-01	1.983	0.3467
14g gin	3.477+00	1.2614-01	0.316	0.7626

Confidentes

Report water is new place new place new place to black the place to be

Ranishal atankeri arror: 13.23 m 377 depans of drawin Rainiple R-opanet: 0.0004,kiperni R-opanet: 0.020 P-stationic: 0.003 m 64 and 377 W, product 0.22016

How about this ????

Here is the Lasso coefficient plot.



Boosting or RF or Deep Learning with SHAP values ???!!!

Still more quesions than answers, *but I would definitely use mBARTD* !!!!!!