Data Mining with Bayesian Trees

Rob McCulloch University of Chicago, Booth School of Business Milwaukee, April 4, 2014

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

Outline:

(i)

Trees and ensemble methods.

(ii)

BART: a Bayesian ensemble method.

(iii)

BART and "big data":

- Parallel version of BART.
- -Consensus Bayes and BART.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

・ロト ・ 日・ ・ 田・ ・ 日・ ・ 日・

Four papers:

Bayesian Additive Regression Trees, Annals of Applied Statistics, 2011. (Chipman, George, McCulloch)

Parallel Bayesian Additive Regression Trees, Journal of Computational and Graphical Statistics, forthcoming.(M. T. Pratola, H. Chipman, J.R. Gattiker, D.M. Higdon, R. McCulloch and W. Rust)

Bayes and Big Data: The Consensus Monte Carlo Algorithm, submitted (Steven L. Scott, Alexander W. Blocker, Fernando V. Bonassi, Hugh A. Chipman, Edward I. George, and Robert E. McCulloch)

Chipman: Acadia; George: U Penn, Wharton; Pratola: Ohio State.

Gattiker, Hidgon, Rust: Los Alamos; Scott, Blocker, Bonassi: Google.

And,

Reversal of fortune: a statistical analysis of penalty calls in the National Hockey League. Journal of Quantitative Analysis in Sports, forthcoming. (Jason Abrevaya, Robert McCulloch)

Abrevaya: University of Texas at Austin.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

The Hockey Data

Glen Healey, commenting on an NHL broadcast:

Referees are predictable. The flames have had three penalties, I guarantee you the oilers will have three.

Well, *guarantee* seems a bit strong, but there is something to it.

How predictable are referees?

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

End

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Got data on every penalty in every (regular season) game for 7 seasons around the time they switched from one referee to two.

For each penalty (after the first one in a game) let

revcall =

1 if current penalty and previous penalty are on different teams,

0 otherwise.

You know a penalty has just been called, which team is it on? is it a reverse call on the other team???

Mean of revcall is .6 !

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

End

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへで

For every penalty (after the first one in a game) we have:

Table: Variable Descriptions

Variable	Description	Mean	Min	Max
Dependent variable				
revcall	1 if current penalty and last penalty are on different teams	0.589	0	1
Indicator-Variable Covariates				
ppgoal	1 if last penalty resulted in a power-play goal	0.157	0	1
home	1 if last penalty was called on the home team	0.483	0	1
inrow2	1 if last two penalties called on the same team	0.354	0	1
inrow3	1 if last three penalties called on the same team	0.107	0	1
inrow4	1 if last four penalties called on the same team	0.027	0	1
tworef	1 if game is officiated by two referees	0.414	0	1
Categorical-variable covariate				
season	Season that game is played		1	7

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

End

◆□▶ ◆□▶ ◆目▶ ◆目▶ 目 のへぐ

Table: Variable Descriptions

Variable	Description	Mean	Min	Max
Other covariates				
timeingame	Time in the game (in minutes)	31.44	0.43	59.98
dayofseason	Number of days since season began	95.95	1	201
numpen	Number of penalties called so far (in the game)	5.76	2	21
timebetpens	Time (in minutes) since the last penalty call	5.96	0.02	55.13
goaldiff	Goals for last penalized team minus goals for opponent	-0.02	-10	10
gf1	Goals/game scored by the last team penalized	2.78	1.84	4.40
ga1	Goals/game allowed by the last team penalized	2.75	1.98	4.44
pf1	Penalties/game committed by the last team penalized	6.01	4.11	8.37
pa1	Penalties/game by opponents of the last team penalized	5.97	4.33	8.25
gf2	Goals/game scored by other team (not just penalized)	2.78	1.84	4.40
ga2	Goals/game allowed by other team	2.78	1.98	4.44
pf2	Penalties/game committed by other team	5.96	4.11	8.37
pa2	Penalties/game by opponents of other team	5.98	4.33	8.25

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

End

n = 57, 883.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへぐ

How is revcall related to the variables?

	inrow2=0	inrow2=1
revcall=0	0.44	0.36
revcall=1	0.56	0.64

inrow2=1:

If the last two calls were on the same team then 64% of the time, the next call will *reverse* and be on the other team.

inrow2=0:

If the last two calls were on different teams, then the frequency of reversal is only 56%.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

Of course,

we want to relate revcall to all the other variables jointly!

Well, we could just run a logit,

but with all the info, can we, fairly automatically, get a better fit than a logit gives?

What could we try?

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

Data Mining Certificates Online Stanford Center for Professional Development

Data Mining and Analysis STATS202 Description

In the Information Age, there is an unprecedented amount of data being collected and stored by banks, supermarkets, internet retailers, security services, etc. So, now that we have all this data, what do we with it?

The discipline of data mining and analysis provides crunchers with the tools and framework to discover meaningful patterns in data sets of any size and scale. It allows us to turn all of this data into valuable, actionable information. In this course, learn how to explore, analyze, and leverage data.

Topics Include

Decision trees Neural networks Association rules Clustering Case-based methods Data visualization

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへで

Modern Applied Statistics: Data Mining STATS315B

Online

Description

Examine new techniques for predictive and descriptive learning using concepts that bridge gaps among statistics, computer science, and artificial intelligence.

This second sequence course emphasizes the statistical application of these areas and integration with standard statistical methodology. The differentiation of predictive and descriptive learning will be examined from varying statistical perspectives.

Topics Include

Classification & regression trees Multivariate adaptive regression splines Prototype & near-neighbor methods Neural networks

Instructors Jerome Friedman, Professor Emeritus, Statistics

Intro

Trees and Ensemble Methods

BAR

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

http://www.sas.com/events/aconf/2010/bdmci61.html

Advanced Analytics for Customer Intelligence Using SAS

Predictive Modeling for Customer Intelligence: The KDD Process Model A Refresher on Data Preprocessing and Data Mining Advanced Sampling Schemes

```
cross-validation (stratified, leave-one-out) bootstrapping
```

Neural networks

multilayer perceptrons (MLPs) MLP types (RBF, recurrent, etc.) weight learning (backpropagation, conjugate gradient, etc.) overfitting, early stopping, and weight regularization architecture selection (grid search, SNC, etc.) input selection (Hinton graphs, likelihood statistics, brute force, etc.) self organizing maps (SOMs) for unsupervised learning case study: SOMs for country corruption analysis

Support Vector Machines (SVMs)

linear programming the kernel trick and Mercer theorem SVMs for classification and regression multiclass SVMs (one versus one, one versus all coding) hyperparameter tuning using cross-validation methods case study: benchmarking SVM classifiers

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

Opening up the Neural Network and SVM Black Box

rule extraction methods (pedagogical versus decompositional approaches such as neurorule, neurolinear, trepan, etc. two-stage models

A Recap of Decision Trees (C4.5, CART, CHAID) Regression Trees

splitting/stopping/assignment criteria

Ensemble Methods

bagging boosting stacking random forests

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ □臣 = のへで

A Tree



- Last penalized was not ahead
- Last two penalties on same team
- Not long since last call
- one ref

72% revcall.

- Last penalized was ahead
- it has been a while since last penalty
- last three calls not on same team

48% revcall.

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

Recall:

How do you fit a tree?

(i)

Build a big tree by greedy search.

(ii)

Prune it back using cross-validation.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

・ロト ・ 日・ ・ 田・ ・ 日・ ・ 日・

Ensemble Methods:

A single tree can be interpretable, but it does not give great in-sample fit or out-of-sample predictive performance.

Ensemble methods combine fit from many trees to give an overall fit.

They can work great!

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

Let's try Random Forests (Leo Brieman).

- Build many (a forest) of very big trees, each of which would overfit on its own.
- Randomize the fitting, so the trees vary (eg. In choosing each decision rule, randomly sample a subset of variables to try, randomly resample data).
- To predict, average (or vote) the result of all the trees in the forest.

For example, build a thousand trees !!!

Have to choose the number of trees in the forest and the randomization scheme.

Wow! (crazy or **brilliant**?)

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

Let's try Random Forests and trees of various sizes.

For Random Forests we have to choose the number of trees to use and the number of variables to sample.

I'll use the default for number of variables and try 200,500, and 1500 trees in the forest.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

・ロト ・ 日・ ・ 田・ ・ 日・ ・ 日・

We have **57,883** observations and a small number of variables so let's do a simple train-test split.

train:

use 47,883 observations to fit the models.

test:

use **10,000** out-of-sample observations to see how well we predict.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

Loss for trees of different sizes, and 3 forests of different sizes (each forest has many big trees!).

Trees and Ensemble Methods

Smaller loss is better.

Loss is measured by the deviance (-2*log-likelihood (out-of-sample)).





We used the default prior and 200 trees.

BART

We want to "fit" the fundamental model:

$$Y_i = f(X_i) + \epsilon_i$$

BART is a Markov Monte Carlo Method that draws from

$$f \mid (x, y)$$

We can then use the draws as our inference for f.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

To get the draws, we will have to:

- Put a prior on f.
- Specify a Markov chain whose stationary distribution is the posterior of f.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

End

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ ● ●

Simulate data from the model:

$$Y_i = x_i^3 + \epsilon_i \quad \epsilon_i \sim N(0, \sigma^2)$$
 iic

n = 100
sigma = .1
f = function(x) {x^3}
set.seed(14)
x = sort(2*runif(n)-1)
y = f(x) + sigma*rnorm(n)
xtest = seq(-1,1,by=.2)

Here, *xtest* will be the *out of sample* \times values at which we wish to infer *f* or make predictions.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes



▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

```
library(BayesTree)
rb = bart(x,y,xtest)
length(xtest)
[1] 11
dim(rb$yhat.test)
[1] 1000 11
```

The (i, j) element of yhat.test is

the i^{th} draw of f evaluated at the j^{th} value of xtest.

1,000 draws of f, each of which is evaluated at 11 xtest values.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

plot(x,y) lines(xtest,xtest^3,col="blue") lines(xtest,apply(rb\$yhat.test,2,mean),col="red") qm = apply(rb\$yhat.test,2,quantile,probs=c(.05,.95)) lines(xtest,qm[1,],col="red",lty=2) lines(xtest,qm[2,],col="red",lty=2)



Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

Example: Out of Sample Prediction

Did out of sample predictive comparisons on 42 data sets. (*thanks to Wei-Yin Loh!!*)

- ▶ p=3-65, n=100-7,000.
- for each data set 20 random splits into 5/6 train and 1/6 test
- use 5-fold cross-validation on train to pick hyperparameters (except BART-default!)
- gives 20*42 = 840 out-of-sample predictions, for each prediction, divide rmse of different methods by the smallest
- + each boxplots represents
 840 predictions for a method
- + 1.2 means you are 20% worse than the best
- + BART-cv best
- BART-default (use default prior) does amazingly well!!



Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

A Regression Tree Model

Let T denote the tree structure including the decision rules.

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

Let g(x; T, M), be a regression tree function that assigns a μ value to x. $x_{5} < C$ $x_{5} \ge C$ $\mu_{3} = 7$ $x_{2} < d$ $x_{2} \ge d$ $\mu_{1} = -2$ $\mu_{2} = 5$

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ ()

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

End

A single tree model:

$$y = g(x; T, M) + \epsilon.$$

A coordinate view of g(x; T, M)



Easy to see that g(x; T, M) is just a step function.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

The BART Model

ntro

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



 $m = 200, 1000, \ldots, \text{big}, \ldots$

 $f(x \mid \cdot)$ is the sum of all the corresponding μ 's at each bottom node.

Such a model combines additive and interaction effects.

Complete the Model with a Regularization Prior

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

 π wants:

- Each T small.
- Each µ small.
- "nice" σ (smaller than least squares estimate).

We refer to π as a regularization prior because it keeps the overall fit from getting "too good".

In addition, it keeps the contribution of each $g(x; T_i, M_i)$ model component small, each component is a "weak learner".

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

BART MCMC

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

First, it is a "simple" Gibbs sampler:

To draw $(T_i, M_i) | \cdot$ we subtract the contributions of the other trees from both sides to get a simple one-tree model. We integrate out M to draw T and then draw M | T. Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

To draw T we use a Metropolis-Hastings within Gibbs step. We use various moves, but the key is a "birth-death" step.



propose a more complex tree

propose a simpler tree

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

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

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



Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Nice things about BART:

- don't have to think about x's (compare: add x_i² and use lasso).
- don't have to prespecify level of interaction (compare: boosting in R)
- competitive out-of-sample.
- stable MCMC.
- stochastic search.
- simple prior.
- uncertainty.
- small p and big n.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

End

・ロト ・ 日本・ 小田・ ・ 田・ うらぐ

How can I use BART to understand the hockey data??

I have a problem !!!.

The sum of trees model is not interpretable.

Other methods (neural nets, random forests, ..) have the same problem.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

・ロト ・ 日・ ・ 田・ ・ 日・ ・ 日・

Well, you can look into the trees for things you can understand:

- what variables get used.
- what variables get used together in the same tree.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Or, estimate p(revcall = 1 | x) at lots of interesting x and see what happens!!

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

In practice, people often make millions of predictions.

We'll keep it simple for the hockey data:

(i)

Change one thing at a time (careful! you can't increase numpens without increasing time-in-game)

(ii)

Do a 2^5 factorial experiment where we move 5 things around, with each thing having two levels. This will give us 32 x vectors to esimate p(revcall = 1 | x) at.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

Change the binary x's one at a time:

Top panel is the posterior distribution of p(revcall) at various x.

Bottom panel is the posterior distribution

of the *difference* in p(revcall) due to a change in x.



Intro

Frees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

▲□▶▲圖▶▲≣▶▲≣▶ ■ のへの

Change the score: *x* is lead of last penalized team.



2⁵ factorial:

r or R:

Last two penalties on different teams or the same team.

g or G:

Last penalized team is behind by one goal or ahead by one goal.

t or T:

Time since last penalty is 2 or 7 minutes.

n or N: Number of penalties is 3 or 12 (time in game is 10 or 55 minutes).

h or H:

Last penalty was not on the home team or it was.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

biggest at gRtnH:

last penalized behind, just had two calls on them, has not been long since last call, early in the game, they are the home team.



Yes, Glen, referees are predictable !!!

If you analyze these carefully, there are interesting interactions!!

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

Top panel: effect of inrow2 at all settings of other variables: *effect is bigger, early in the game.*



Bottom panel: effect of goaldiff at all settings of other variables

effect is bigger, late in the game.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

PBART: Parallel Bayesian Additive Trees

Dave Higdon said,

we tried your stuff (the R package BayesTree) on the analysis of computer experiments and it seemed promising but it is too slow".

Recode with MPI to make it faster!!

MPI: Message Passing Interface.

Two Steps.

Step 1. Rewrote serial code so that it is "leaner".

Intro

Trees and Ensemble Methods

BAR

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

End

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・ うへぐ

Lean Code:

```
class tree {
public:
. . .
private:
  //-----
  //parameter for node
  double mu;
  //-----
  //rule: left if x[v] < xinfo[v][c]</pre>
  size_t v;
  size t c:
  //-----
  //tree structure
  tree_p p; //parent
  tree_p l; //left child
  tree_p r; //right child
};
```

1 double, two integers, three pointers.

ntro

Trees and Ensemble Methods

BAR

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

n	bart/BayesTree MCMC	new MCMC
1,000	57.725	14.057
2,000	136.081	27.459
3,000	211.799	40.483
4,000	298.712	54.454
5,000	374.971	66.900
6,000	463.861	82.084
7,000	545.995	95.737
8,000	651.683	107.911
9,000	724.577	120.778
10,000	817.711	135.764

ntro

Frees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

End

The new code is 4 to 6 times faster!!

> 57.725/14.057
[1] 4.106495
> 817.711/135.764
[1] 6.023033

Note: We also use a more limited set of MCMC moves in the new code but we find we get the same fits.

Step 2:

Parallel MPI implementation.

Have p + 1 processor cores. Split data up into p equal chunks.

- core 0 is the master. It runs the MCMC.
- core *i*, i = 1, 2, ..., p has data chunk *i* in memory.
- ► Each core has the complete model ((T_j, M_j)^m_{j=1}, σ) in memory.

Note: with MPI cores and associated memory may be on different machines.

Compare with openmp where the memory must be shared.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

- core 0 is the master. It runs the MCMC.
- core i, i = 1, 2, ..., p has data chunk i in memory.
- Each core has the complete model $((T_j, M_j)_{j=1}^m, \sigma)$ in memory.

Each MCMC step involves:

- 1. master core 0, initiates an MCMC step (e.g. change a single tree, draw σ).
- master core 0, sends out a compute request (needed for MCMC step) to each slave node i = 1, 2, ... p.
- 3. Each slave core *i* computes on it's part of the data and sends the results back to master core 0.
- 4. master core 0 combines the results from the *p* slaves and updates the model using the results (e.g. changes a tree, obtains new σ draw).
- 5. master core 0, copies new model state out to all the slave cores.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

Keys to Parallel implementation:

- Lean model representation is cheap to copy out to all the slave cores.
- MCMC draws all depend on simple conditionally sufficient statistics which may be computed on the separate slave cores, cheaply sent back to the master, and then combined.

Even though the the overall model is complex, each local move is simple !!

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

End

・ロト ・ 日・ ・ 田・ ・ 日・ ・ 日・

Simple Sufficient Statistics:

Consider the birth-death step:



Given a tree, we just have $\{R_{ij}\} \sim N(\mu_j, \sigma^2)$ iid in j^{th} bottom node, where R is resids from the the other trees.

Evaluating a birth-death is just like testing equality of two normal means with an independent normal prior.

Sufficient statistic is just $\sum_{i} r_{ij}$ for two different *j* corresponding to left/right child bottom nodes.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

Timing

We did lots of timing exercises to see how well it works.

Here n = 200,000, p + 1 = 40, time is to do 20,000 iterations.

processors	run time (s)
2	347087
4	123802
8	37656
16	16502
24	9660
30	6303
40	4985
48	4477

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

End

40/8=5. 39/7 = 5.6. 37656/4985=7.5.

With 5 times as many processors we are 7.5 times faster.

Here is a plot of the times.

 T_{ser} : time with serial code; T_{par} : time with parallel code. Ideally you might hope for

7

$$\overline{p_{par}} = \frac{T_{ser}}{p+1}.$$

So we plot
$$\frac{1}{T_{par}}$$
 vs. $p+1$.

Looks pretty linear.





Intro

Trees and Ensemble Methods

BAR

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

Here we look at the efficiency which is

$${\sf E} = rac{{{T_{ser}}}}{{\left({p + 1}
ight){T_{par}}}}$$

If $T_{par} = \frac{T_{ser}}{(p+1)}$ this would equal 1.

We timed 27 runs with $m \in \{50, 100, 200\}$, $n \in \{100000, 500000, 1000000\}$, $p + 1 \in \{9, 17, 25\}$.

If you have too few observations on a core the cost of message passing eats into the speed.

60% of p + 1 times faster is still pretty good!



Intro

Trees and Ensemble Methods

BAR

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

Consensus Bayes

At Google they were interested, but also politely scoffed.

Their attitude is that on a typical large network the communication costs will really bite.

Some of our timings were on a shared memory machine with 32 cores over which we had complete control, but some were on a cluster where you buy computes and these are over several machines.

In any case, our setup may not work with the kind of machine setup they want to use at Google.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

So, a simpler idea is consensus Bayes.

Suppose you have a "cluster" of p "nodes".

Again, you split the data up into p chunks (no master this time!).

Then you assign each node one of the chunks of data.

Then you simply run a separate MCMC on each "node" using its chunk of data.

If you want to predict at x, let $f_{ij}(x)$ be the i^{th} draw of f on node j.

Then you get a set of consensus draws by

$$f_i^c(x) = \sum w_j f_{ij}(x), \quad w_j = 1/\mathsf{Var}(f_{ij}(x)).$$

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

End

◆□▶ ◆□▶ ◆目▶ ◆目▶ 目 のへぐ

Consensus Bayes also involves an adjustment of the prior. The idea is that

$$p(\theta \mid D) \propto L(D)p(\theta) \propto \prod_{j=1}^{p} L(D_j)p(\theta)^{1/p}.$$

where D is all the data and D_i is the data on chunk j.

So, you should use prior $p(\theta)^{1/p}$ on each node.

This can work really well for simple parametric models, but BART is "non-parametric" with variable dimension.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

Simulated data.

n = 20,000. k = 10. p = 30.

- True f(x) vs. serial BART f(x).
 True f(x) vs. consensus BART f(x), no prior adjustment.
- (3) True f(x) vs. consensus BART $\hat{f}(x)$, prior adjustment.



Consensus BART without prior adjustment is awesome!

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

◆□ ▶ ◆□ ▶ ◆ 三 ▶ ◆ 三 ● ● ● ●

Top row: consensus no prior adjustment. Bottom row: consensus prior adjustment.

First column, posterior mean of f(x) vs. serial BART posterior mean. Second column, 5% quantile vs serial BART. Third column, 95% quantile vs serial BART.



Consensus BART without prior adjustment is awesome!

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Frees

Consensus Bayes

End

Good:

Both approaches make BART a usable technique with large data sets and seem to work very well.

Need to think more about how to make prior adjustment (if any) for Consensus Bayes.

Bad:

While PBART is on my webpage, we need to make it easier to use.

We have some R functions written and are testing, but it is not available yet.

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

Very Bad:

Have been promising new version of package for a while.

New faster version (with parallel computing) will be there soon!!

Intro

Trees and Ensemble Methods

BART

PBART: Parallel Bayesian Additive Trees

Consensus Bayes

End

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ ● ●