

Package: stan4bart (via r-universe)

September 6, 2024

Version 0.0-10

Date 2024-09-03

Title Bayesian Additive Regression Trees with Stan-Sampled Parametric Extensions

Depends R (>= 3.5-0), methods, dbarts (>= 0.9-21)

Imports stats, Matrix, parallel, RcppParallel (>= 5.1.1)

LinkingTo BH (>= 1.72.0.3), Rcpp (>= 1.0.5), RcppEigen (>= 0.3.3.7.0), RcppParallel (>= 5.1.1), dbarts (>= 0.9-20)

Suggests testthat (>= 2.0-0), lme4

Description Fits semiparametric linear and multilevel models with non-parametric additive Bayesian additive regression tree (BART; Chipman, George, and McCulloch (2010) <[doi:10.1214/09-AOAS285](https://doi.org/10.1214/09-AOAS285)>) components and Stan (Stan Development Team (2021) <<https://mc-stan.org/>>) sampled parametric ones. Multilevel models can be expressed using 'lme4' syntax (Bates, Maechler, Bolker, and Walker (2015) <[doi:10.18637/jss.v067.i01](https://doi.org/10.18637/jss.v067.i01)>).

License GPL (>= 3)

NeedsCompilation true

Biarch true

UseLTO true

URL <https://github.com/vdorie/stan4bart>

BugReports <https://github.com/vdorie/stan4bart/issues>

Repository <https://vdorie.r-universe.dev>

RemoteUrl <https://github.com/vdorie/stan4bart>

RemoteRef HEAD

RemoteSha 42f9c5857604b99e19cf38d1c896c87b973424f6

Contents

stan4bart	2
stan4bart-generics	8

Index	11
--------------	-----------

stan4bart	<i>Semiparametric Models Using Stan and BART</i>
-----------	--

Description

This function fits semi-parametric linear and probit models that have a non-parametric, BART component and one or more of a parametric fixed effect (unmodeled coefficients), or a parametric random effect (modeled coefficients). If $f(x)$ is a BART “sum-of-trees” [model](#), fits:

- For continuous response variables:

$$Y | b \sim N(f(X^b) + X^f \beta + Zb, \sigma^2) \quad b \sim N(0, \Sigma_b)$$

- For binary response variables:

$$P(Y = 1 | b) = \Phi(f(X^b) + X^f \beta + Zb) \quad b \sim N(0, \Sigma_b)$$

Usage

```
stan4bart(formula,
  data = NULL,
  subset,
  weights,
  na.action = getOption("na.action", "na.omit"),
  offset,
  contrasts = NULL,
  test = NULL,
  treatment = NULL,
  offset_test = NULL,
  verbose = FALSE,
  iter = 2000L,
  warmup = iter %% 2L,
  skip = 1L,
  chains = 4L,
  cores = getOption("mc.cores", 1L),
  refresh = max(iter %% 10L, 1L),
  offset_type = c("default", "fixef", "ranef", "bart", "parametric"),
  seed = NA_integer_,
  keep_fits = TRUE,
  callback = NULL,
  stan_args = NULL,
  bart_args = NULL)
```

Arguments

formula	a formula object, or one that can be coerced to that type. Terms on the right-hand-side of the formula that are encased in a symbolic call to <code>bart()</code> will be used to create the non-parametric component. Terms that use the lmer -style grouping syntax will be added as parametric, hierarchical varying intercepts and slopes. All other terms will be added as fixed effects.
data	an optional data frame containing the variables in the formula. Its use is strongly encouraged.
subset, weights, na.action, offset, contrasts	optional components adjusting the constructed model matrices and otherwise changing the linear predictor. <code>na.action</code> cannot be <code>"na.pass"</code> . See lm and model.matrix.default .
test	an optional data frame to be used as test data. If present, the test predictions will be stored as the sampler runs and can be extracted later.
treatment	an optional symbol, that when present and refers to a binary variable, will be used to create a test data frame with the treatment variable set to its counterfactual. Only one of <code>test</code> and <code>treatment</code> can be supplied.
offset_test	optional vector which will be added to the test predictions.
verbose	a logical or integer. If <code>FALSE</code> or non-positive, runs quietly. Additional levels of information may be displayed for increasingly positive numbers, however a large number of diagnostics are suppressed when running multi-threaded. If negative, all diagnostic information is ignored.
iter	positive integer indicating the number of posterior samples to draw and return. Includes warmup samples.
warmup	non-negative integer indicating number of posterior samples to draw and throw away. Also controls the number of iterations for which adaptation is run.
skip	one or two positive integers. Every <code>skip</code> sample will be kept, while every other sample will be discarded. If argument is length two, an attempt will be made to use the named element <code>"bart"</code> for BART and <code>"stan"</code> for Stan. If not named, BART is the first skip element and Stan is the second. This argument does not impact the number of <code>iters</code> returned, unlike a conventional "thinning" parameter.
chains	positive integer giving the number of Markov Chains to sample.
cores	positive integer giving the number of units of parallelization. Computation for each chain will be divide among the cores available. When greater than one, verbose output within chains will not be available.
refresh	positive integer giving the frequency with which diagnostic information should be printed as the sampler runs. Only applies with <code>cores</code> (or <code>chains</code>) equal to 1.
offset_type	character; an experimental/testing feature that controls how <code>offset</code> is to be interpreted. When one of <code>"fixef"</code> , <code>"ranef"</code> , or <code>"bart"</code> , the offset is used to replace that part of the model. When <code>"parametric"</code> , it replaces both of the fixed and random parametric components. Sampling is still done for these components and their draws are stored, however whenever they were present in the fit the supplied value is used instead.

seed	Optional integer specifying the desired pRNG seed . It should not be needed when running single-threaded - calling <code>set.seed</code> will suffice. The primary use of seed is to obtain reproducible results when multi-threaded. See Reproducibility section below.
keep_fits	Logical that, when false, prevents the sampler from storing each draw. Intended to be used with <code>callback</code> .
callback	A function that will be called at each iteration, accepting three arguments: <code>yhat.train</code> , <code>yhat.test</code> , <code>stan_pars</code> . See details for more information.
stan_args	optional list, specifying further arguments to Stan. See details below.
bart_args	optional list, specifying further arguments to BART. See details below.

Details

Fits a Bayesian “mixed effect” model with a non-parametric Bayesian Additive Regression Trees (BART) component. For continuous responses:

- $Y_i | b \sim N\left(f(X_i^b) + X_i^f \beta + Z_i b_{g[i]}, \sigma^2\right)$
- $b_j \sim N(0, \Sigma_b)$

where b_j are the “random effects” - random intercepts and slopes - that correspond to group j , $g[i]$ is a mapping from individual i to its group index, f - a BART sum-of-trees model, X^b are predictors used in the BART model, X^f are predictors in a parametric, linear “fixed effect” component, Z is the design matrix for the random intercept and slopes, and σ and Σ_b are variance components.

Binary outcome models are obtained by assuming a latent variable that has the above distribution, and that the observed response is 1 when that variable is positive and 0 otherwise. The response variable marginally has the distribution:

- $$P(Y_i = 1 | b) = \Phi\left(f(X_i^b) + X_i^f \beta + Z_i b_{g[i]}\right)$$

where Φ is the cumulative distribution function of the standard normal distribution.

Terminology: As `stan4bart` fits a Bayesian model, essentially all components are “modeled”. Furthermore, as it has two first-level, non-hierarchical components, “fixed” effects are ambiguous. Thus we adopt:

- “fixed” - refers only to the parametric, linear, individual level mean component, $X^f \beta$; these are “unmodeled coefficients” in other contexts
- “random” - refers only to the parametric, linear, hierarchical mean component, Zb ; these are “modeled coefficients” in other contexts
- “bart” - refers only to the nonparametric, individual level mean component, $f(X^b)$

Model Specification: Model specification occurs in the `formula` object, according to the following rules:

- variables or terms specified inside a pseudo-call to `bart` are used for the “bart” component, e.g. `y ~ bart(x_1 + x_2)`

- variables or terms specified according to [lmer](#) syntax are used for the “random” effect component, e.g. $y \sim (1 \mid g_1) + (1 + x_3 \mid g_1)$
- remaining variables not inside a `bart` or “bars” construct are used for the “fixed” effect component; e.g. $y \sim x_4$

All three components can be present in a single model, however a `bart` part must be present. If you wish to fit a model without one, use `stan_glmmer` in the `rstanarm` package instead.

Additional Arguments: The `stan_args` and `bart_args` arguments to `stan4bart` can be used to pass further arguments to `stan` and `bart` respectively. These are similar to the functions `stan` in the `rstan` package and `bart`, but not identical as `stan4bart` constructs its own model internally. Stan arguments include:

- `prior_covariance`
- `prior`, `prior_intercept`, `prior_aux`, `QR`
- `init_r`, `adapt_gamma`, `adapt_delta`, `adapt_kappa` - see the help page for `stan` in the `rstan` package.

For reference on the first two sets of options, see the help page for `stan_glmmer` in the `rstanarm` package; for reference on the third set, see the help page for `stan` in the `rstan` package.

BART arguments include:

- further arguments to `dbartsControl` that are not specified by `stan4bart`, such as `keepTrees` or `n.trees`; keeping trees can be costly in terms of memory, but is required to use `predict`

Reproducibility: Behavior differs when running multi- and single-threaded, as the pseudo random number generators (pRNG) used by R are not thread safe. When single-threaded, R’s built-in generator is used; if set at the start, `.Random.seed` will be used and its value updated as samples are drawn. When multi-threaded, the default behavior is draw new random seeds for each thread using the clock and use thread-specific pRNGs.

This behavior can be modified by setting `seed`. For the single-threaded case, that seed will be installed and the existing seed replaced at the end, if applicable. For multi-threaded runs, the seeds for threads are drawn sequentially using the supplied seed, and will not change the state of R’s built-in generator.

Consequently, the `seed` argument should not be needed when running single-threaded - `set.seed` will suffice. When multi-threaded, `seed` can be used to obtain reproducible results.

Callbacks: Callbacks can be used to avoid expensive memory allocation, aggregating results as the sampler proceeds instead of waiting until the end. A callback function must accept the arguments:

- `yhat.train` - BART predictions of the expected value of the training data, conditioned on the Stan parameters
- `yhat.test` - when applicable, the same values as above but for the test data; NULL otherwise
- `stan_pars` - named vector of Stan samples, including fixed and random effects and variance parameters

It is expected that the callback will return a vector of the same length each time. If not, invalid memory writes will likely result. If the result of the callback has names, they will be added to the result.

Serialization: At present, `stan4bart` models cannot be safely `saved` and `loaded` in a way that the sampler can be restarted. This feature may be added in the future.

Value

Returns a list assigned class `stan4bartFit`. Has components below, some of which will be NULL if not applicable.

Input values:

<code>y</code>	response vector
<code>weights</code>	weights vector or null
<code>offset</code>	offset vector or null
<code>frame</code>	joint model frame for all components
<code>formula</code>	formula used to specify the model
<code>na.action</code>	supplied <code>na.action</code>
<code>call</code>	original call

Stored data:

<code>bartData</code>	data object used for BART component
<code>X</code>	fixed effect design matrix or NULL
<code>X_means</code>	column means of fixed effect design matrix when appropriate
<code>reTrms</code>	random effect “terms” object when applicable, as used by lmer
<code>test</code>	named list when applicable, having components <code>X</code> and <code>reTrms</code> ; test data for BART is added to the <code>bartData</code> result
<code>treatment</code>	treatment vector, when applicable

Results, better accessed using [extract](#):

<code>bart_train</code>	samples of individual posterior predictions for BART component
<code>bart_test</code>	predicted test values for BART component, when applicable
<code>bart_varcount</code>	BART variable counts
<code>sigma</code>	samples of residual standard error; not present for binary outcomes
<code>k</code>	samples of the end-node sensitivity parameter; only present when it is modeled
<code>ranef</code>	samples of random effects, or modeled coefficients; will be a named list, with effects for each grouping factor
<code>Sigma</code>	samples of covariance of random effects; also a named list with one element for each grouping factor
<code>fixef</code>	samples of the fixed effects, or unmodeled coefficients
<code>callback</code>	samples returned by an optional callback function

Other items:

<code>warmup</code>	a list of warmup samples, containing the same objects in the results subsection
<code>diagnostics</code>	Stan sampler produced diagnostic information, include tree depth and divergent transitions
<code>sampler.bart</code>	external points to BART samplers; used only for predict when <code>keepTrees</code> is TRUE
<code>range.bart</code>	internal scale used by BART samplers, used by predict when <code>keepTrees</code> is TRUE

Author(s)

Vincent Dorie: <vdorie@gmail.com>.

See Also

[bart](#), [lmer](#), and [stan_glm](#) in the [rstanarm](#) package

Examples

```
# simulate data (extension of Friedman MARS paper)
# x consists of 10 variables, only first 5 matter
# x_4 is linear
f <- function(x)
  10 * sin(pi * x[,1] * x[,2]) + 20 * (x[,3] - 0.5)^2 +
  10 * x[,4] + 5 * x[,5]

set.seed(99)
sigma <- 1.0

n <- 100
n.g.1 <- 5L
n.g.2 <- 8L

# sample observation level covariates and calculate marginal mean
x <- matrix(runif(n * 10), n, 10)
mu.bart <- f(x) - 10 * x[,4]
mu.fixef <- 10 * x[,4]

# varying intercepts and slopes for first grouping factor
g.1 <- sample(n.g.1, n, replace = TRUE)
Sigma.b.1 <- matrix(c(1.5^2, .2, .2, 1^2), 2)
b.1 <- matrix(rnorm(2 * n.g.1), n.g.1) %%% chol(Sigma.b.1)

# varying intercepts for second grouping factor
g.2 <- sample(n.g.2, n, replace = TRUE)
Sigma.b.2 <- as.matrix(1.2)
b.2 <- rnorm(n.g.2, 0, sqrt(Sigma.b.2))

mu.ranef <- b.1[g.1,1] + x[,4] * b.1[g.1,2] + b.2[g.2]

y <- mu.bart + mu.fixef + mu.ranef + rnorm(n, 0, sigma)

df <- data.frame(y, x, g.1, g.2)

fit <- stan4bart(
  formula = y ~
    X4 + # linear component ("fixef")
    (1 + X4 | g.1) + (1 | g.2) + # multilevel ("ranef")
    bart(. - g.1 - g.2 - X4), # use bart for other variables
  verbose = -1, # suppress ALL output
  # low numbers for illustration
  data = df,
```

```

chains = 1, iter = 10, bart_args = list(n.trees = 5))

# posterior means of individual expected values
y.hat <- fitted(fit)

# posterior means of the random effects
ranef.hat <- fitted(fit, type = "ranef")

```

stan4bart-generics *Generic Functions for stan4bart Model Fits*

Description

Commonly expected utility functions to derive useful quantities from fitted models.

Usage

```

## S3 method for class 'stan4bartFit'
extract(
  object,
  type = c("ev", "ppd", "fixef", "indiv.fixef", "ranef", "indiv.ranef",
           "indiv.bart", "sigma", "Sigma", "k", "varcount", "stan",
           "trees", "callback"),
  sample = c("train", "test"),
  combine_chains = TRUE,
  sample_new_levels = TRUE,
  include_warmup = FALSE,
  ...)

## S3 method for class 'stan4bartFit'
fitted(
  object,
  type = c("ev", "ppd", "fixef", "indiv.fixef", "ranef", "indiv.ranef",
           "indiv.bart", "sigma", "Sigma", "k", "varcount", "stan",
           "callback"),
  sample = c("train", "test"),
  sample_new_levels = TRUE,
  ...)

## S3 method for class 'stan4bartFit'
predict(
  object, newdata, offset,
  type = c("ev", "ppd", "indiv.fixef", "indiv.ranef", "indiv.bart"),
  combine_chains = TRUE,
  sample_new_levels = TRUE,
  ...)

```


Arguments

object	a fitted model resulting from a call to <code>stan4bart</code> .
type	a character vector; one of the options listed below.
sample	one of "train" or "test", indicating if the training or test data frames should be used.
combine_chains	logical controlling if chain information should be discarded and the result returned as a matrix instead of an array.
sample_new_levels	logical; if TRUE, levels out of the training sample will have random effects drawn from their posterior predictive distribution. If FALSE, their random effects will be fixed to 0.
include_warmup	logical or "only"; when TRUE/FALSE, warmup samples will or will not be included in the result respectively. When "only", only the warmup samples will be returned.
newdata	data frame for making out of sample predictions.
offset	optional vector which will be added to test predictors.
...	not currently in use, but provided to match signatures of other generics.

Details

`extract` is used to obtain raw samples using the training or test data, `fitted` averages those samples, and `predict` operates on data not available at the time of fitting. Note: `predict` requires that the model be fit with `args_bart = list(keepTrees = TRUE)`.

Return type: The type argument accepts:

- "ev" - the individual level expected value, that is draws from $E[Y \mid X^b, X^f, Z] \mid Y = f(X^b) + X^f\beta + Zb \mid Y$ where the expectation is with respect to the posterior distribution of the parameters given the data
- "ppd" - draws from the individual level posterior predictive distribution, generally speaking adding noise to the result for "ev" or simulating new Bernoulli trials.
- "fixef" - draws from the posterior of the fixed effects (also known as the "unmodeled" coefficients), $\beta \mid Y$
- "indiv.fixef" - draws from the posterior distribution of the individual level mean component deriving from the fixed effects, $X^f\beta$
- "ranef" - the random effects, varying intercepts and slopes, or "modeled" coefficients, b ; b has substantial structure that is represented as the returned value, where coefficients are reported within their grouping factors
- "indiv.ranef" - individual level mean component deriving from the random effects, Zb
- "indiv.bart" - individual level mean component deriving from the BART model, $f(X^b)$
- "sigma" - for continuous responses, the residual standard error
- "Sigma" - when applicable, the covariance matrices of the random effects
- "stan" - raw matrix or array of Stan sampled transformed parameters.
- "trees" - a data frame of flattened trees; see the subsection on extracted trees in [bart](#) and note that `stan4bart` variable names can be found in the `bartData@x` element of a fitted `stan4bart` model

- "callback" - if a callback function was provided while fitting, the results of that for each sample

Value

extract and predict return either arrays of dimensions equal to $n.\text{observations} \times n.\text{samples} \times n.\text{chains}$ when `combine_chains` is `FALSE`, or matrices of dimensions equal to $n.\text{observations} \times (n.\text{samples} * n.\text{chains})$ when `combine_chains` is `TRUE`.

fitted returns a vector of the appropriate length by averaging the result of a call to extract.

Author(s)

Vincent Dorie: <vdorie@gmail.com>.

Index

.Random.seed, 5
bart, 5, 7, 9
data, 6
dbartsControl, 5
extract, 6
extract (stan4bart-generics), 8
fitted.stan4bartFit
 (stan4bart-generics), 8
formula, 3
lm, 3
lmer, 3, 5–7
load, 5
model, 2
model.matrix.default, 3
predict, 5, 6
predict.stan4bartFit
 (stan4bart-generics), 8
save, 5
seed, 4
set.seed, 4, 5
stan4bart, 2, 9
stan4bart-generics, 8