

APPLICATION

embarcadero: Species distribution modelling with Bayesian additive regression trees in R

Colin J. Carlson^{1,2} ¹Department of Biology, Georgetown University, Washington, DC, USA²Center for Global Health Science and Security, Georgetown University Medical Center, Washington, DC, USA**Correspondence**

Colin J. Carlson

Email: colin.carlson@georgetown.edu

Handling Editor: Samantha Price**Abstract**

1. *embarcadero* is an R package of convenience tools for species distribution modelling (SDM) with Bayesian additive regression trees (BART), a powerful machine learning approach that has been rarely applied to ecological problems.
2. Like other classification and regression tree methods, BART estimates the probability of a binary outcome based on a set of decision trees. Unlike other methods, BART iteratively generates sets of trees based on a set of priors about tree structure and nodes, and builds a posterior distribution of estimated classification probabilities. So far, BARTs have yet to be applied to SDM.
3. *embarcadero* is a workflow wrapper for BART species distribution models, and includes functionality for easy spatial prediction, an automated variable selection procedure, several types of partial dependence visualization and other tools for ecological application. The *embarcadero* package is an open source and available on Github.
4. To show how *embarcadero* can be used by ecologists, I illustrate a BART workflow for a virtual species distribution model. The supplement includes a more advanced vignette showing how BART can be used for mapping disease transmission risk, using the example of Crimean–Congo haemorrhagic fever in Africa.

KEYWORDS

Bayesian additive regression trees, Crimean–Congo haemorrhagic fever, disease ecology, ecological niche modelling, machine learning, population ecology, regression trees, species distribution modelling

1 | INTRODUCTION

In the last two decades, over two dozen statistical and machine learning methods have been applied to species distribution modelling (SDM; Norberg et al., 2019). Over time, a handful of methods have risen to predominance due to ease of implementation, computational speed and strong predictive performance in rigorous cross-validation. Some methods are especially popular for specific applications, mostly because of disciplinary tradition. For example, maximum entropy (MaxEnt) models are widely popular for studies of global ecological responses to climate change (VanDerWal et al.,

2013; Warren et al., 2013). In disease ecology, boosted regression trees (BRTs) have become the dominant tool for mapping vectors, reservoirs and transmission risk of infectious zoonoses and vector-borne diseases (Carlson et al., 2019; Messina et al., 2016; Pigott et al., 2014), largely due to an influential 2013 paper on dengue virus (Bhatt et al., 2013). SDMs are used for several—sometimes conflicting—purposes in ecology, and popular methods are sometimes used despite known shortcomings (Guillera-Aroita et al., 2015; Smith & Santos, 2019). In particular, most popular methods have a limited framework for handling uncertainty, and conspicuously few popular methods are Bayesian (and vice versa).

In this paper, I discuss a new Bayesian approach to classification and regression trees (CART), one of the most popular families of machine learning methods used in ecology. Models in this family estimate the probability of a given output variable (in this case, a binary classification of habitat suitability or species presence) based on decision 'trees' that split predictor variables with nested, binary rule-sets. The precise rules for generating these trees vary across implementations. For example, in *random forest* models, an ensemble of trees is generated, where each tree is independently generated based on a bootstrap of the original dataset; trees grow to the maximum possible depth (the longest chain of splitting rules), with no pruning (trees are never post hoc reduced). In the BRTs approach, shallower trees with a constrained depth ('weak learners') are constructed iteratively that explain the residuals left by previous trees; this adds bias, but allows the model to focus on unusual cases at the potential expense of overfitting (Elith, Leathwick, & Hastie, 2008; Vezhnevets & Barinova, 2007). CART methods have many strengths for SDM; they consistently perform well in model comparisons (Elith et al., 2006; Mainali et al., 2015; Redding, Lucas, Blackburn, & Jones, 2017; Wisz et al., 2008), and the tree-based approach is often more intuitive than the complex fitting procedures 'under the hood' of MaxEnt or Maxlike methods (Elith et al., 2011; Merow & Silander, 2014; Merow, Smith, & Silander, 2013).

Bayesian additive regression trees (BART) are an exciting and new alternative to other popular classification tree methods. As in other approaches, BART generates a set of decision trees that explain different components of variance in the outcome variable. Unlike random forests or BRTs, the formulation of BART is Bayesian, with the posterior probability of a model shaped by priors $P(\text{trees})$ on how trees should look (i.e. the parameters used to generate those trees):

$$P(\text{trees}|\text{data}) \propto P(\text{data}|\text{trees}) P(\text{trees}). \quad (1)$$

Like BRTs, BART introduces variance by fitting a set of many shallow 'weak learner' trees, but unlike BRT, this is explicitly controlled by three prior distributions: the probability a tree stops at a node of a given depth, the probability of a given variable being drawn for a splitting rule and the probability of splitting that variable at a particular value. The latter two are usually treated as uniformly distributed (splits are randomly constructed by variable, and within each variables' range), while the first is usually specified as a negative power law, constraining tree depth and penalizing overfitting. Using these priors, a specified number of trees m are generated with no splits, and then updated randomly in an MCMC process that allows them to be expanded, rearranged or pruned. Each model instance is a *sum-of-trees* model, unlike random forests, which average predictions across trees; to create the sum-of-trees model, each tree is adjusted to the residuals of the sum-of-remaining-trees. This process superficially resembles how boosting works within BRTs, but because trees are tuned to the ensemble, they rarely overfit to particular cases within the residuals (Chipman, George, & McCulloch, 2010). After a fixed number of burn-in samples that are dropped, the full set of sum-of-trees models across all samples from the Markov

chain is treated as a posterior distribution and used to generate the posterior distribution of predictions (For a more in-depth explanation, including a visualization of tree structure in the MCMC process, see Tan & Roy, 2019).

In computer science, BARTs are used for everything from medical diagnostics to self-driving car algorithms (Sparapani et al., 2018; Tan, Flannagan, & Elliott, 2018); however, they have yet to find widespread application in ecology. A study from 2011 used BART as a tool to examine habitat selection data on birds (Yen, Thomson, Vesk, & Mac Nally, 2011); a 2017 study used BART to evaluate performance data of other SDM methods (Farley, 2017). But so far, they have not been used for the purpose of predicting species distributions. This reflects a broader deficit of Bayesian models in the SDM literature: several elegant Bayesian SDM methods have been previously proposed (Golding & Purse, 2016; Redding et al., 2017), but none are particularly widely adopted, possibly because advanced Bayesian models may seem discouraging or unintuitive.

Bayesian additive regression trees brings the conceptual familiarity and strengths of classification tree methods, but adds a relatively simple Bayesian component that inherently and intuitively handles model uncertainty. This might make it a promising alternative not just to existing Bayesian approaches but also popular classification tree methods, in particular BRTs. BRT has several easy to use out of the box implementations, is powerful for ecological inference and consistently performs well in rigorous tests of SDM performance. However, BRT also has downsides: it can be prone to overfitting, and fitting procedures are largely handed down as anecdotal best practices, with many studies choosing hyperparameters based on software defaults; very few studies select parameters from formal cross-validation as early work recommended (Elith et al., 2008). Furthermore, uncertainty is usually measured by generating an unweighted ensemble of BRT submodels over sub-setted training data, generating a confidence interval from data permutations (like random forests) rather than formal assumptions about model uncertainty. In contrast, the formal Bayesian structure of BART captures uncertainty within a single model, which is more coherent and intuitive than how uncertainty is usually generated in BRT ensembles. BART also shares many of the strengths of BRT, like easy out of the box implementation and easy visualization of 'black box' model components, and outperforms other CART methods in model comparisons (Chipman et al., 2010).

This paper introduces an R package, *embarcadero*, as a convenience tool for running SDMs with BARTs. Throughout, I use a simulated 'virtual species' (see Appendix 1 in Supporting Information) to illustrate the workflow and the major features of the package, including model selection, visualization and diagnostics. Because BRTs are the most popular method of SDM in medical geography, the supplement includes a second, more detailed vignette using BART to map Crimean–Congo haemorrhagic fever (CCHF) in Africa, based on the distribution of the tick *Hyalomma truncatum*, a presumed vector. This is a more challenging and computationally intensive implementation, and takes several hours to run on most machines, but

highlights some of the strength of the approach for applied scientific questions.

2 | SDMs WITH BARTs

2.1 | Implementing BART with binary classification

At least four R packages currently exist that can implement BARTs: *BayesTree* (Chipman & McCulloch, 2016), *bartMachine* (Kapelner & Bleich, 2013), *BART* (McCulloch, Sparapani, Gramacy, Spanbauer, & Pratola, 2018) and *dbarts* (Chipman, McCulloch, & Dorie, 2014). Their functionality differs in important ways, and not all of them are currently capable of important features like partial dependence plots that are important for SDMs. This package is an SDM-oriented workflow wrapper for *dbarts*, which includes most of the basic functionality needed for SDM, including a simple implementation of BART with binary outcomes. A list of the functions made available in *embarcadero*, versus their counterparts and additional useful functions in *dbarts*, is given in Table 1.

In the original notation of Chipman et al. (2010), BART consists of tree structures T and terminal nodes (leaves) M , as an ensemble $(T_1, M_1), \dots, (T_m, M_m)$. Each tree generates a predictive function $g(\cdot)$, with a sum-of-trees function $f(\cdot)$ given as

$$f(\cdot) = \sum_{j=1}^m g(\cdot; T_j, M_j) + \epsilon; \quad \epsilon \sim \mathcal{N}(0, \sigma^2). \quad (2)$$

A set of posterior draws f^* of sum-of-trees functions, generated by the MCMC process described above, create the posterior distribution for $p(f|y) \equiv P(\text{trees}|\text{data})$. Given the assumption of normality in the outcome variable and error structure, BART handles binary classification problems (like SDM) using a logit link, where Φ is the standard normal c.d.f. and:

$$f(\cdot) = \Phi \left[\sum_{j=1}^m g(\cdot; T_j, M_j) \right]. \quad (3)$$

Binary classification is run by `dbarts::bart` automatically when supplied with a binary outcome. However, the returned predictions are untransformed back into probabilities, a problem solved in *embarcadero* with a `predict` wrapper (This also allows prediction on raster datasets, a key piece of SDM workflow).

2.2 | An example of a BART SDM

To see how BART works, we can generate a virtual species on a hypothetical landscape which responds to climate variables X1 through X4, but is uninfluenced by variables X5 to X8 (see Appendix 1 in Supporting Information). Like most other SDM methods in R, the BART model itself is run on a data frame of presence-absence or presence-pseudoabsence points, and associated environmental

TABLE 1 Functions available in *embarcadero* and additional functions in *dbarts* of importance

Core modelling functionality	
<code>dbarts::bart</code>	Runs a binary Bayesian additive regression trees (BART) classification model
<code>bart.step</code>	Full implementation of a BART model with built-in variable set reduction (a wrapper for <code>dbarts::bart</code> , <code>variable.step</code> , <code>varimp</code> , <code>varimp.diag</code> , and <code>summary</code>)
<code>predict</code>	Predict species distributions with a BART model and a <i>RasterStack</i> of environmental layers (a wrapper for <code>dbarts::predict.bart</code>)
<code>summary</code>	Returns a summary of call, performance and diagnostic plots for a BART model object
<code>dbarts::xbart</code>	Cross-validation of hyperparameters, in particular for the tree depth prior, with several metrics of model performance
Variable diagnostics	
<code>variable.step</code>	Stepwise variable set reduction algorithm
<code>varimp</code>	Returns variable importance, with optional plots
<code>varimp.diag</code>	Diagnostic of variable importance at different m values
Visualization	
<code>partial</code>	Partial dependence plots for single variables (a <i>ggplot2</i> -based wrapper for <code>dbarts::pdbart</code>)
<code>dbarts::pd2bart</code>	Two-predictor, three-dimensional partial dependence plots (no wrapper implemented yet)
<code>plot.mcmc</code>	Visualize each posterior draw's prediction and the running average of those predictions. Can be used with the <i>animation</i> package to create GIFs of how the posterior draw learns to fit the data (especially interesting for the burn-in of models with small number of trees)
<code>spartial</code>	Spartial projection (maps) of partial dependence plots onto raw environmental covariates
Convenience tools	
<code>bigstack</code>	Fast aggregation of an environmental layer <i>RasterStack</i> for quick prediction, using the <i>velox</i> package

covariates. For example, with a *RasterStack* called `climate` and an occurrence dataset called `occ.df`, the basic workflow is

```
library(embarcadero)
xnames <- c("x1", "x2", "x3", "x4",
            "x5", "x6", "x7", "x8")

## Run the BART model
sdm <- bart(y.train=occ.df[, "Observed"],
            x.train=occ.df[, xnames],
            keeptrees = TRUE)
```

```
## Predict the species distribution
map <- predict(sdm, climate)
## Visualize model performance
summary(bart)
```

This last line returns a brief model diagnostic including the optimal cut-off for thresholding classifications and some measures of performance, like the area under the receiver-operator curve (AUC):

```
Call: bart occ.df[, xnames] occ.df[, "Observed"] TRUE
Predictor list:
  x1 x2 x3 x4 x5 x6 x7 x8
Area under the receiver-operator curve
  AUC = 0.91
Recommended threshold (maximizes true skill statistic)
  Cutoff = 0.42
  TSS = 0.71
Resulting type I error rate: 0.078
Resulting type II error rate: 0.21
```

Additionally, `summary` returns a diagnostic figure (Figure 1), summarizing the performance of the classifier on the training data.

The primary appeal of BART, compared to other CART methods, is a formal way of measuring model uncertainty within any individual implementation. Visualizing uncertainty in BART predictions is easy

with `embarcadero`; for example, to derive a 95% credible interval from the 2.5% and 97.5% quantiles of the posterior, a user can specify:

```
map <- predict(sdm, climate, quantiles=c(0.025, 0.975))
```

Instead of returning a single raster, the function returns a stack of rasters with each specified quantile. Mapping the difference between the two quantile rasters gives the credible interval width, which provides a native measure of spatial uncertainty, analogous to how the coefficient of variation can be used to measure spatial uncertainty across an ensemble of BRT runs (Carlson et al., 2019). When running tasks especially with several quantiles, or large rasters, prediction runtime grows quickly and memory can become limiting; `predict` has a `splitby` option that breaks the task into pieces, which minimizes memory conflicts, adds a progress bar and allows estimation of total runtime based on the first chunk:

```
map <- predict(sdm, climate, quantiles=c(0.025,
0.975), splitby=10)
```

3 | VARIABLE SELECTION

Variable importance (calculated by `varimp`) is usually measured in BART models by counting the number of times a given variable is used by a tree split across the full posterior draw of trees (This

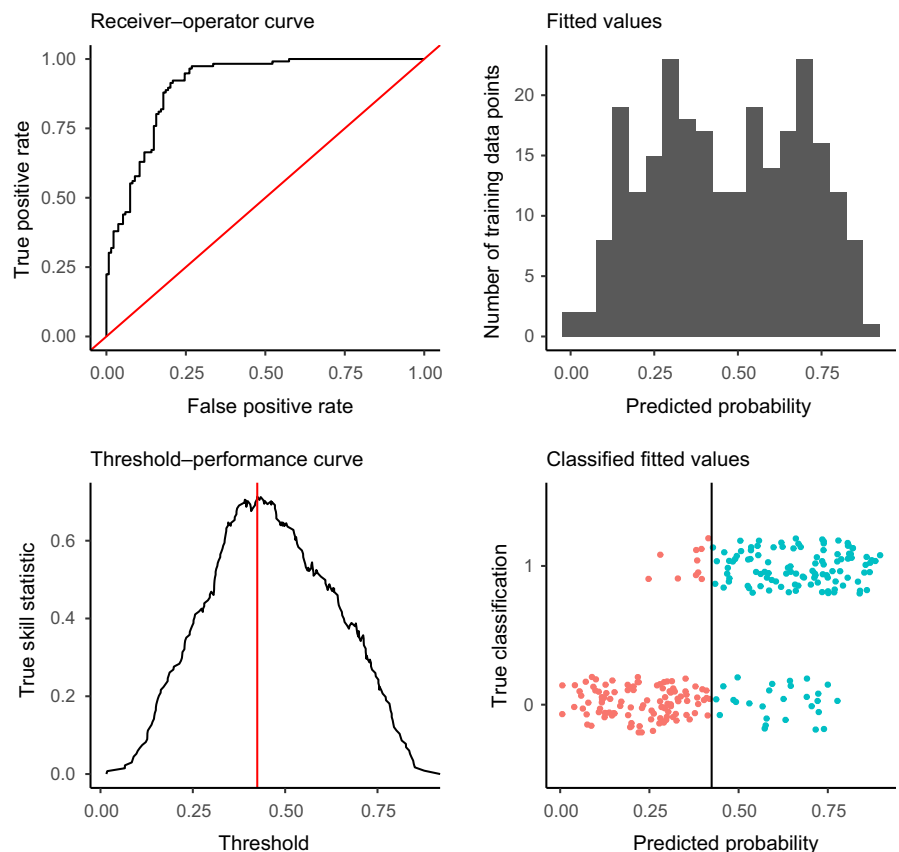


FIGURE 1 The model diagnostic returned by `summary`. A high area under the receiver-operator curve and clear visual split in predicted probabilities assigned to the training presences and pseudoabsences indicates that the model has done an adequate job

is similar to variable importance in BRTs, which is calculated from the number of tree splits and the corresponding improvement they cause in the model). In models with higher numbers of trees, the difference in variable importance becomes less pronounced, and less informative variables receive a higher number of splitting rules. Conversely, variable selection can be performed by running models with a small number of trees ($m = 10$ or 20), and observing which variables stop being included in trees (Chipman et al., 2010). This diagnostic is generated in *embarcadero* by `varimp.diag` (see an example in Figure 2).

Analysis of this diagnostic plot is still subjective and informal. As a way to standardize variable set reduction rules across workflows, *embarcadero* includes an automatic variable selection procedure in `variable.step`:

1. Fit a full model with all predictors and a small tree ensemble (default $m = 10$), a fixed number of times (default $n = 50$);
2. Eliminate the least informative variable across all 50 runs;
3. Rerun the models again minus the least informative variable ($n = 50$ times again), recording the root mean square error (RMSE; on the training data);
4. Repeat steps 2 and 3 until there are only three covariates left;
5. Finally, select the model with the lowest average RMSE.

Anecdotally, this procedure almost always recommends dropping every variable with decreasing importance in models with fewer trees, and conserves every variable with increasing importance. In our virtual species case, for example, the diagnostic shows that X1 through X4 have much higher performance than X5 through X8 (Figure 2), and the automated procedure recommends dropping X5 through X8:

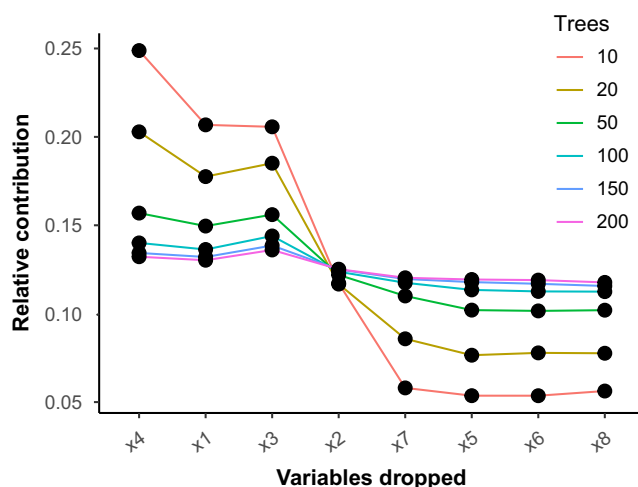


FIGURE 2 The model diagnostic returned by `varimp.diag`. The four informative predictors (X1 through X4) stay in the model as the number of trees m is reduced down to 10; the four uninformative predictors (X5 through X8) contribute to the models only minimally, or are entirely dropped. The diagnostic therefore allows a user to correctly identify the four environmental covariates that actually drive the species' distribution

```
varimp.diag(occ.df[,xnames],
            occ.df[, "Observed"],
            iter=50)

step.model <- variable.step(x.data=occ.df[,xnames],
                           y.data=occ.df[, "Observed"])

step.model
[1] "x1" "x2" "x3" "x4"
```

This largely matches original work which found that BART is highly effective at identifying informative subsets of predictors (see section 5.2.1 of Chipman et al., 2010).

I recommend careful analysis of all diagnostic information, but include a full automated variable selection pipeline in `bart.step`, which (a) produces the initial multi- m diagnostic plot, (b) runs automated variable selection, (c) returns a model trained with the optimal variable set, (d) plots variable importance in the final model and (e) returns the summary of the final model. Despite automation, this procedure is not a fail-safe against the inclusion of uninformative predictors, or false inference on them; this is true of almost all methods, and predictors should always be chosen based on at least some expert opinion about biological plausibility (Fourcade, Besnard, & Secondi, 2018). Similarly, validation of partial dependence curves against biological knowledge should be treated as an additional level of model validation, potentially more informative than measuring predictive accuracy (Warren, Matzke, & Iglesias, 2019).

4 | VISUALIZING MODEL RESULTS

embarcadero includes several methods for generating partial dependence plots. The function `partial` is written as a wrapper for `dbarts::pdpbart`, and can be used to generate partial dependence plots with a customizable, `ggplot2`-based aesthetic, including multiple ways of visualizing uncertainty (As with overall predictions, credible intervals on partial plots are true Bayesian credible intervals). Posteriors can be visualized with traceplots of individual draws, or bars for a credible interval of a specified width (by default 95%):

```
partial(sdm, x.vars=c("x4"),
       smooth=5,
       equal=TRUE,
       trace=FALSE)

## VERSUS, for comparison,
gbml <- dismo::gbm.step(data=occ.df,
                       gbm.x = 2:5, gbm.y = 1,
                       family = "bernoulli",
                       tree.complexity = 5,
                       learning.rate = 0.01,
                       bag.fraction = 0.5)

dismo::gbm.plot(gbml, variable.no=4, rug=TRUE,
                plot.layout=c(1,1))
```


This visualizes uncertainty much clearer than, for example, `dismo::gbm.plot` can in a single instance (Figure 3). Two-dimensional partial dependence plots (interactions among two predictor variables) can also be generated using `dbarts::pd2bart`.

Finally, `embarcadero` includes a new visualization called *spatial partial dependence plots*, which reclassify predictor rasters based on their partial dependence plots, and show the relative suitability of different regions for an individual covariate. The `spartial` function can be used to generate these maps, and answer questions like 'What desert regions are too arid, even in their wettest month, for spade-foot toads'? or 'Where are the soils with the best pH for redwood

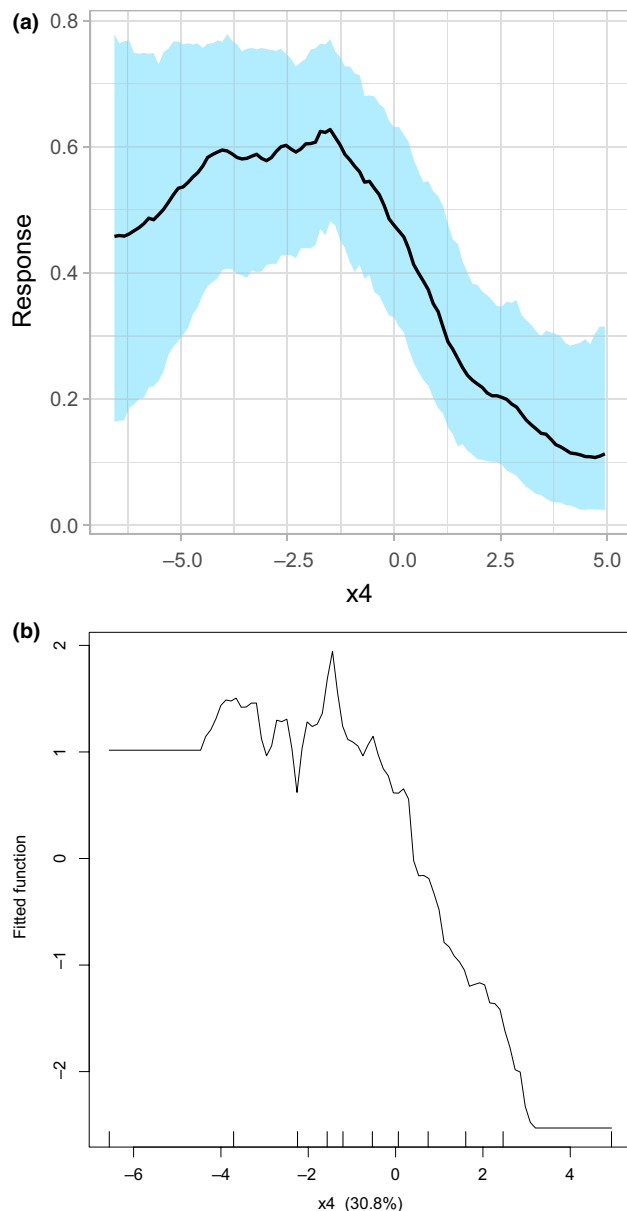


FIGURE 3 Partial dependence curves generated by single-instance Bayesian additive regression trees implementations (a) show uncertainty with more transparency and clarity than those generated from single-instance BRT implementations (b)

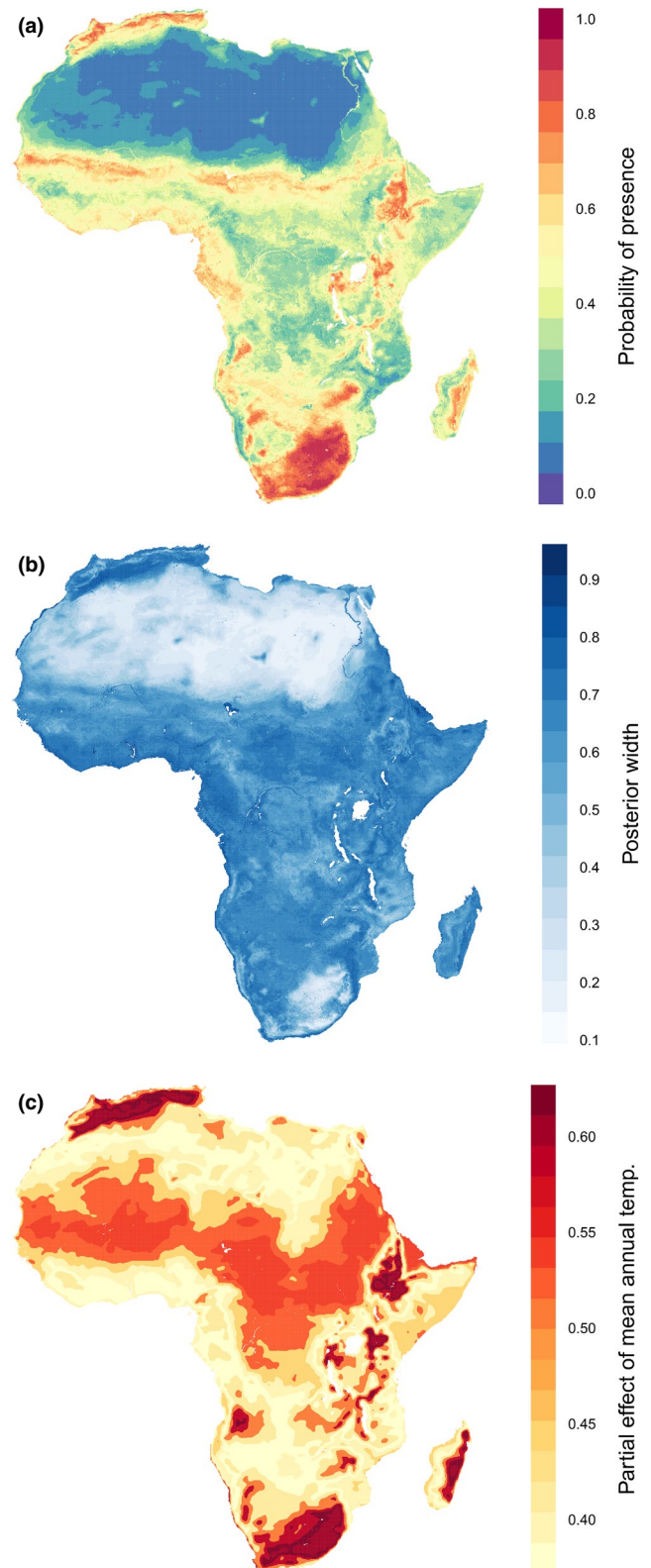


FIGURE 4 A map of Crimean-Congo haemorrhagic fever transmission risk, constructed using ecological niche modelling with Bayesian additive regression trees (see Supporting Information), including (a) the posterior mean, (b) the posterior width (95% credible interval) and (c) a spatial dependence plot for BIO1 (mean annual temperature)

growth'? These visualization options are illustrated in greater depth in the advanced vignette (see Figure 4c).

5 | AN ADVANCED VIGNETTE

To demonstrate applications to disease transmission mapping, the supplement includes an advanced tutorial on *embarcadero* focused on updating an African risk map for CCHF virus, is a tick-borne Bunyavirus that causes extremely severe, and often fatal, illness in humans. Very little is known about CCHF, compared to other cosmopolitan tick-borne illnesses like Lyme disease or tularemia. The definitive reservoir of CCHF is unknown but likely ungulates (Babayan, Orton, & Streicker, 2018); outbreaks frequently affect sheep and other domestic ruminants. The vectors of CCHF are better known, and are presumed to almost always be *Hyalomma* ticks, which are widespread throughout Africa and Eurasia; other tick vectors have been suspected, but evidence for their competence is limited (Papa, Tsergouli, Tsioka, & Mirazimi, 2017). In Africa, *H. truncatum* in particular is common throughout rangeland and is a strong candidate for a primary vector (Logan, Linthicum, Bailey, Watts, & Moulton, 1989; Wilson, Gonzalez, Cornet, & Camicas, 1991). A global map of CCHF has been previously produced with BRTs; a significant amount of the Black Sea region was suitable, while areas outside had highly localized predictions of suitability, presumably because of data sparsity in Africa especially (Messina, Pigott, Golding, et al., 2015). However, some major areas of presence appeared under-predicted, such as the western Congo Basin.

The advanced vignette shows how BART can be used to map CCHF in Africa, using the same occurrence dataset as previous mapping efforts have (Messina, Pigott, Duda, et al., 2015). Just as studies of dengue risk have included suitability for the *Aedes aegypti* mosquito as a covariate, the new model includes a suitability layer for *H. truncatum*, created from the canonical dataset on African tick distributions (Cumming, 1998). The updated map predicts that the distribution of CCHF may be more geographically expansive than previous studies have indicated (Figure 4). Areas of the highest risk are still heavily concentrated in Sahel rangeland and east African highlands, but also far more extensive in southern Africa and along the Atlantic coast than previously believed. A detailed tutorial is provided showing this workflow in the Supporting Information of this paper, and all data are available online (github.com/cjcarlson/pier39).

6 | DISCUSSION

Because BART is a comparatively new method, many of the basic use case questions remain mostly unaddressed: Do pseudoabsences perform notably worse than absences? Is there a minimum sample size? Does collinearity inflate or distort variable importance? Users may wish to explore some of these points using virtual species before working with BART on their data, or to compare BART results to other methods as a sense check.

Furthermore, as with any other Bayesian method, out of the box implementation can make it easy to neglect or underconsider prior selection. Beginning users can use `dbarts::xbart` for a cross-validation-based approach to hyperparameter selection for the tree depth prior. More advanced users may be interested in going more in depth within the BART literature to set better priors. For example, using a uniform prior on covariate importance can be unhelpful—especially in high-dimensionality data with only a few valid predictors, where the model tends to converge on the variable importance prior (Rockova & van der Pas, 2017; Tan et al., 2018). Instead, setting a Dirichlet distribution for the prior (DART) can significantly improve model performance and variable selection (Linero, 2018). These models have also been extended for smooth response functions (softBART and softDART), and often serve as a challenge for CART models, which may be particularly relevant to the types of concave response functions expected for ecophysiological traits (Linero & Yang, 2018).

Finally, it is worth mentioning that BART is a growing topic of interest in machine learning, and new extensions may expand applications within SDM work and more broadly in spatial ecology. For example, the random intercept BART (riBART) model is a framework for handling cases of structure within outcome data; this framework might be useful for cases where sampling bias has categorical structure (e.g. different levels of sampling across country or state borders; Tan et al., 2018). Preliminary compatibility exists in the package for prediction using the `dbarts::rbart_vi` function. Similarly, causal inference using the BART framework has become especially popular (Hahn, Murray, & Carvalho, 2017), which may be an interesting direction for modelling given recent work proposing causal inference as a new priority for mapping infectious diseases (Kraemer, Reiner, & Bhatt, 2019). Expanding work along these lines will help establish better best practices for using BARTs in SDM applications.

ACKNOWLEDGEMENTS

Thanks to Vincent Dorie for creating the fabulous and wonderfully useful *dbarts* package, to Shweta Bansal and Jason Blackburn for research support, to Ethan Beaman for telling me about BARTs in the first place, to Zack Susswein for helpful comments on design of Bayesian models, to Tad Dallas for helpful comments on the R package and manuscript and to David Lawrence Miller and a second anonymous reviewer for detailed and helpful feedback. This work was funded by a Georgetown Environment Initiative (GEI) postdoctoral fellowship.

DATA AVAILABILITY STATEMENT

No original data are presented in the paper, which includes code for data simulation. In the Supplement, a series of datasets are used with original citations given and the raw files available in the vignette repository at github.com/cjcarlson/pier39 archived on Zenodo: <https://doi.org/10.5281/zenodo.3703234> (Carlson, 2020a). Version 1.1 of the R package itself is also archived on Zenodo: <https://doi.org/10.5281/zenodo.3703228> (Carlson, 2020b).

ORCID

Colin J. Carlson  <https://orcid.org/0000-0001-6960-8434>

REFERENCES

- Babayan, S. A., Orton, R. J., & Streicker, D. G. (2018). Predicting reservoir hosts and arthropod vectors from evolutionary signatures in RNA virus genomes. *Science*, 362, 577–580. <https://doi.org/10.1126/science.aap9072>
- Bhatt, S., Gething, P. W., Brady, O. J., Messina, J. P., Farlow, A. W., Moyes, C. L., ... Hay, S. I. (2013). The global distribution and burden of dengue. *Nature*, 496, 504–507. <https://doi.org/10.1038/nature12060>
- Carlson, C. J. (2020a). Data from: cjcarlson/pier39: ME&E Paper release - advanced vignette (Version v1.0). *Zenodo*, <https://doi.org/10.5281/zenodo.3703234>
- Carlson, C. J. (2020b). Data from: cjcarlson/embarcadero: Publication version (ME&E) (Version v1.1.0.0001). *Zenodo*, <http://doi.org/10.5281/zenodo.3703228>
- Carlson, C. J., Kracalik, I. T., Ross, N., Alexander, K. A., Hugh-Jones, M. E., Fegan, M., ... Blackburn, J. K. (2019). The global distribution of *Bacillus anthracis* and associated anthrax risk to humans, livestock and wildlife. *Nature Microbiology*, 4, 1337–1343. <https://doi.org/10.1038/s41564-019-0435-4>
- Chipman, H. A., George, E. I., & McCulloch, R. E. (2010). BART: Bayesian additive regression trees. *The Annals of Applied Statistics*, 4, 266–298. <https://doi.org/10.1214/09-AOAS285>
- Chipman, H., & McCulloch, R. (2016). *BayesTree: Bayesian additive regression trees*. R package version 0.3-1.3. Retrieved from <https://cran.r-project.org/web/packages/BayesTree/index.html>
- Chipman, H., McCulloch, R., & Dorie, V. (2014). *dbarts: Discrete Bayesian additive regression trees sampler*. R package version 0.8-5. Retrieved from <https://cran.r-project.org/web/packages/dbarts/index.html>
- Cumming, G. (1998). Host preference in African ticks (Acari: Ixodida): A quantitative data set. *Bulletin of Entomological Research*, 88, 379–406. <https://doi.org/10.1017/S0007485300042139>
- Elith, J., H. Graham, C., P. Anderson, R., Dudík, M., Ferrier, S., Guisan, A., ... E. Zimmermann, N. (2006). Novel methods improve prediction of species' distributions from occurrence data. *Ecography*, 29, 129–151. <https://doi.org/10.1111/j.2006.0906-7590.04596.x>
- Elith, J., Leathwick, J. R., & Hastie, T. (2008). A working guide to boosted regression trees. *Journal of Animal Ecology*, 77, 802–813. <https://doi.org/10.1111/j.1365-2656.2008.01390.x>
- Elith, J., Phillips, S. J., Hastie, T., Dudík, M., Chee, Y. E., & Yates, C. J. (2011). A statistical explanation of MaxEnt for ecologists. *Diversity and Distributions*, 17, 43–57. <https://doi.org/10.1111/j.1472-4642.2010.00725.x>
- Farley, S. S. (2017). *A general framework for predicting the optimal computing configurations for climate-driven ecological forecasting models* (PhD thesis). Madison, WI: University of Wisconsin.
- Fourcade, Y., Besnard, A. G., & Secondi, J. (2018). Paintings predict the distribution of species, or the challenge of selecting environmental predictors and evaluation statistics. *Global Ecology and Biogeography*, 27, 245–256. <https://doi.org/10.1111/geb.12684>
- Golding, N., & Purse, B. V. (2016). Fast and flexible Bayesian species distribution modelling using Gaussian processes. *Methods in Ecology and Evolution*, 7, 598–608. <https://doi.org/10.1111/2041-210X.12523>
- Guillera-Aroita, G., Lahoz-Monfort, J. J., Elith, J., Gordon, A., Kujala, H., Lentini, P. E., ... Wintle, B. A. (2015). Is my species distribution model fit for purpose? Matching data and models to applications. *Global Ecology and Biogeography*, 24, 276–292. <https://doi.org/10.1111/geb.12268>
- Hahn, P. R., Murray, J. S., & Carvalho, C. (2017). Bayesian regression tree models for causal inference: Regularization, confounding, and heterogeneous effects. *arXiv preprint arXiv:1706.09523*.
- Kapelnr, A., & Bleich, J. (2013). *bartMachine: Machine learning with Bayesian additive regression trees*. *arXiv preprint arXiv:1312.2171*.
- Kraemer, M. U., Reiner Jr., R. C., & Bhatt, S. (2019). Causal inference in spatial mapping. *Trends in Parasitology*, 35, 743–746. <https://doi.org/10.1016/j.pt.2019.06.005>
- Linero, A. R. (2018). Bayesian regression trees for high-dimensional prediction and variable selection. *Journal of the American Statistical Association*, 113, 626–636. <https://doi.org/10.1080/01621459.2016.1264957>
- Linero, A. R., & Yang, Y. (2018). Bayesian regression tree ensembles that adapt to smoothness and sparsity. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 80, 1087–1110. <https://doi.org/10.1111/rssb.12293>
- Logan, T. M., Linthicum, K. J., Bailey, C. L., Watts, D. M., & Moulton, J. R. (1989). Experimental transmission of Crimean–Congo hemorrhagic fever virus by *Hyalomma truncatum* Koch. *The American Journal of Tropical Medicine and Hygiene*, 40, 207–212. <https://doi.org/10.4269/ajtmh.1989.40.207>
- Mainali, K. P., Warren, D. L., Dhileepan, K., McConnachie, A., Strathie, L., Hassan, G., ... Parmesan, C. (2015). Projecting future expansion of invasive species: Comparing and improving methodologies for species distribution modeling. *Global Change Biology*, 21, 4464–4480. <https://doi.org/10.1111/gcb.13038>
- McCulloch, R., Sparapani, R., Gramacy, R., Spanbauer, C., & Pratola, M. (2018). *BART: Bayesian additive regression trees*. R package version 1.0. Retrieved from <https://cran.r-project.org/web/packages/BART/index.html>
- Merow, C., & Silander, J. A. (2014). A comparison of maxlike and maxent for modelling species distributions. *Methods in Ecology and Evolution*, 5, 215–225. <https://doi.org/10.1111/2041-210X.12152>
- Merow, C., Smith, M. J., & Silander, J. A. (2013). A practical guide to maxent for modeling species' distributions: What it does, and why inputs and settings matter. *Ecography*, 36, 1058–1069. <https://doi.org/10.1111/j.1600-0587.2013.07872.x>
- Messina, J. P., Kraemer, M. U. G., Brady, O. J., Pigott, D. M., Shearer, F. M., Weiss, D. J., ... Hay, S. I. (2016). Mapping global environmental suitability for Zika virus. *eLife*, 5, e15272. <https://doi.org/10.7554/eLife.15272>
- Messina, J. P., Pigott, D. M., Duda, K. A., Brownstein, J. S., Myers, M. F., George, D. B., & Hay, S. I. (2015). A global compendium of human Crimean–Congo haemorrhagic fever virus occurrence. *Scientific Data*, 2, 150016. <https://doi.org/10.1038/sdata.2015.16>
- Messina, J. P., Pigott, D. M., Golding, N., Duda, K. A., Brownstein, J. S., Weiss, D. J., ... Hay, S. I. (2015). The global distribution of Crimean–Congo hemorrhagic fever. *Transactions of the Royal Society of Tropical Medicine and Hygiene*, 109, 503–513. <https://doi.org/10.1093/trstmh/trv050>
- Norberg, A., Abrego, N., Blanchet, F. G., Adler, F. R., Anderson, B. J., Anttila, J., ... Ovaskainen, O. (2019). A comprehensive evaluation of predictive performance of 33 species distribution models at species and community levels. *Ecological Monographs*, 89(3), e01370. <https://doi.org/10.1002/ecm.1370>
- Papa, A., Tsergouli, K., Tsioka, K., & Mirazimi, A. (2017). Crimean–Congo hemorrhagic fever: Tick-host-virus interactions. *Frontiers in Cellular and Infection Microbiology*, 7, 213. <https://doi.org/10.3389/fcimb.2017.00213>
- Pigott, D. M., Golding, N., Mylne, A., Huang, Z., Henry, A. J., Weiss, D. J., ... Hay, S. I. (2014). Mapping the zoonotic niche of Ebola virus disease in Africa. *eLife*, 3, e04395. <https://doi.org/10.7554/eLife.04395>
- Redding, D. W., Lucas, T. C., Blackburn, T. M., & Jones, K. E. (2017). Evaluating Bayesian spatial methods for modelling species distributions with clumped and restricted occurrence data. *PLoS ONE*, 12, e0187602. <https://doi.org/10.1371/journal.pone.0187602>
- Rockova, V., & van der Pas, S. (2017). Posterior Concentration for Bayesian Regression Trees and Forests. *arXiv:1708.08734*. Retrieved from <https://arxiv.org/abs/1708.08734>
- Smith, A. B., & Santos, M. J. (2019). Testing the ability of species distribution models to infer variable importance. *bioRxiv*, 715904.
- Sparapani, R., Dabbouseh, N., Gutterman, D., Zhang, J., Chen, H., Bluemke, D., ... Soliman, E. (2018). Novel electrocardiographic

- criteria for the diagnosis of left ventricular hypertrophy derived with Bayesian additive regression trees: The multi-ethnic study of atherosclerosis. *Circulation*, 138, A10908.
- Tan, Y. V., Flannagan, C. A., & Elliott, M. R. (2018). Predicting human-driving behaviour to help driverless vehicles drive: Random intercept Bayesian additive regression trees. *Statistics and Its Interface*, 11, 557–572. <https://doi.org/10.4310/sii.2018.v11.n4.a1>
- Tan, Y. V., & Roy, J. (2019). Bayesian additive regression trees and the general BART model. arXiv preprint arXiv:190107504.
- VanDerWal, J., Murphy, H. T., Kutt, A. S., Perkins, G. C., Bateman, B. L., Perry, J. J., & Reside, A. E. (2013). Focus on poleward shifts in species' distribution underestimates the fingerprint of climate change. *Nature Climate Change*, 3, 239–243. <https://doi.org/10.1038/nclimate1688>
- Vezhnevets, A., & Barinova, O. (2007). Avoiding boosting overfitting by removing confusing samples. In J. N. Kok, J. Koronacki, R. Lopez de Mantaras, S. Matwin, D. Mladenić, & A. Skowron (Eds.), *European Conference on Machine Learning* (pp. 430–441). Berlin, Germany: Springer.
- Warren, D. L., Matzke, N. J., & Iglesias, T. L. (2019). Evaluating species distribution models with discrimination accuracy is uninformative for many applications. *BioRxiv*, 684399.
- Warren, R., VanDerWal, J., Price, J., Welbergen, J. A., Atkinson, I., Ramirez-Villegas, J., ... Lowe, J. (2013). Quantifying the benefit of early climate change mitigation in avoiding biodiversity loss. *Nature Climate Change*, 3, 678. <https://doi.org/10.1038/nclimate1887>
- Wilson, M., Gonzalez, J. P., Cornet, J. P., & Camicas, J. L. (1991). Transmission of Crimean–Congo haemorrhagic fever virus from experimentally infected sheep to *Hyalomma truncatum* ticks. *Research in Virology*, 142, 395–404. [https://doi.org/10.1016/0923-2516\(91\)90007-P](https://doi.org/10.1016/0923-2516(91)90007-P)
- Wisz, M. S., Hijmans, R. J., Li, J., Peterson, A. T., Graham, C., Guisan, A., & NCEAS Predicting Species Distributions Working Group. (2008). Effects of sample size on the performance of species distribution models. *Diversity and Distributions*, 14, 763–773. <https://doi.org/10.1111/j.1472-4642.2008.00482.x>
- Yen, J. D., Thomson, J. R., Vesk, P. A., & Mac Nally, R. (2011). To what are woodland birds responding? Inference on relative importance of in-site habitat variables using several ensemble habitat modelling techniques. *Ecography*, 34, 946–954. <https://doi.org/10.1111/j.1600-0587.2011.06651.x>

SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section.

How to cite this article: Carlson CJ. *embarcadero*:

Species distribution modelling with Bayesian additive regression trees in R. *Methods Ecol Evol*. 2020;11:850–858.

<https://doi.org/10.1111/2041-210X.13389>