

The forestecology R package for modeling interspecies competition between trees

Albert Y. Kim *

Program in Statistical & Data Sciences, Smith College

and

David Allen

Biology Department, Middlebury College

and

Simon P. Couch

Mathematics Department, Reed College

February 12, 2021

Abstract

Move abstract below here after completed.

Keywords: forest ecology, competition, R, Rstats, tidyverse, sf, cross-validation,

*Albert Y. Kim is Assistant Professor, Statistical & Data Sciences, Smith College, Northampton, MA 01063 (e-mail: akim04@smith.edu).

¹ **Abstract (350 words)**

- ² 1. When modeling growth of trees forest ecologists often incorporate the effect of in-
³ terspecies competition. Many such models are based on a neighborhood effect as-
⁴ sumption whereby all trees within a fixed distance of all focal trees are considered
⁵ competitors. Methods are needed to evaluate the effect of interspecies competition
⁶ and to assess their quality.
- ⁷ 2. We present the `forestecology` package providing methods for both 1) evaluating the
⁸ out-of-sample performance of our model using spatial-crossvalidation and 2) testing a
⁹ null hypothesis of no impact of competitor species' identity on the growth of trees
¹⁰ using a permutation test. We implement a class and methods using R's S3 object-
¹¹ oriented system, for a specific linear, Bayesian neighborhood competition model of
¹² tree growth.
- ¹³ 3. We demonstrate the package's functions using data from the Smithsonian Conserva-
¹⁴ tion Biology Institute's large forest dynamics plot, part of the ForestGEO network
¹⁵ of research sites. Given ForestGEO's data collection protocols and data formatting
¹⁶ standards, the package cross-compatibility of code. We show both that 1) competi-
¹⁷ tor species identity matters and 2) that not spatially cross-validating leads to error
¹⁸ estimates that are overly optimistic.
- ¹⁹ 4. The package follows `tidyverse`-like structure whereby verb-named functions can be
²⁰ modularly “piped” in sequence to intuitively display the sequence of steps of analysis
²¹ from start to finish. Additionally, most inputs/outputs of functions assume an are
²² of `sf` class from the simple features package, thereby facilitating all wrangling and
²³ visualization of geospatial data. Lastly, even though our package is currently limited
²⁴ to one specific model, the package is setup such that it can be easily extended to
²⁵ other models.

²⁶ **1 Introduction**

²⁷ Repeat-censused forest plots offer excellent data to test neighborhood models of tree com-
²⁸ petition Allen & Kim (2020) Canham et al. (2006) Uriarte et al. (2004). Here we describe

29 an R package, `forestecology`, to do that. This package implements the methods in Allen
30 & Kim (2020). It provides: a convenient way to specify and fit models of tree growth based
31 on neighborhood competition; a spatial cross validation method to test and compare model
32 fits Roberts et al. (2017); and an ANOVA-like method to assess whether the competitor
33 identity matters in these models. The model is written to work with ForestGEO plot data
34 Anderson-Teixeira et al. (2015), but we envision that it could easily be modified to work
35 with data from other forest plots, e.g. the US Forest Service Forest Inventory and Analysis
36 plots Smith (2002).

37 The `forestecology` is designed with “tidy” data principles in mind as Wickham et al.
38 (2019).

39 Given that our data is of geo-spatial nature, we represent our data using the “simple
40 features” `sf` package class of objects Pebesma (2018) whereby. While previously the `sp`
41 package serves such purposes Pebesma & Bivand (2005), the `sf` package is designed to
42 interface with the `tidyverse` suite of packages.

43 2 Example

44 We demonstrate the `forestecology` package’s features on the Smithsonian Conservation
45 Biology Institute (SCBI) large forest dynamics plot, located at the Smithsonian’s National
46 Zoo and Conservation Biology Institute in Front Royal, VA, USA. The 25.6 ha (640 x 400
47 m) plot is located at the intersection of three of the major physiographic provinces of the
48 eastern US: the Blue Ridge, Ridge and Valley, and Piedmont provinces and is adjacent to
49 the northern end of Shenandoah National Park. The forest type is typical mature secondary
50 eastern mixed deciduous forest, with a canopy dominated by tulip poplar (*Liriodendron*
51 *tulipifera*), oaks (*Quercus* spp.), and hickories (*Carya* spp.), and an understory composed
52 mainly of spicebush (*Lindera benzoin*), paw-paw (*Asimina triloba*), American hornbeam
53 (*Carpinus caroliniana*), and witch hazel (*Hamamelis virginiana*) Bourg et al. (2013).

54 A high-level overview of the steps of our analysis pipeline is as follows:

- 55 1. Compute the growth of trees based on census data
- 56 2. Add spatial information:

- 57 1. Define buffer region trees.
58 2. Add spatial cross-validation block information.

59 3. Compute focal versus competitor tree information.
60 4. Fit model and make predictions.
61 5. Additionally: Evaluate model performance using spatial cross-validation.
62 6. Additionally: Evaluate the effect of competitor species identity using permutation
63 tests.

64 We load all necessary packages.

```
library(tidyverse)
library(lubridate)
library(sf)
library(forestecology)
```

65 **2.1 Compute the growth of trees based on census data**

66 The first step in the our analysis sequence is to compute the growth of trees using data
67 from two censuses. The `compute_growth()` function computes growth assuming census
68 data that follows ForestGEO standards. Despite such standards, minor variations will still
69 exist between sites thereby necessitating some data wrangling and checking. For example,
70 the SCBI site records all dbh's in millimeters, whereas the Michigan Big Woods site records
71 them in centimeters Anderson-Teixeira et al. (2015) Allen et al. (2020). The data format of
72 other sites may be such that our `compute_growth()` function doesn't work at all. However,
73 in the end all that matters is that the growth of all trees is saved in a data frame of type
74 class `sf` whereby the geolocation of each tree is presented in a `geometry` variable of type
75 `<POINT>` and at a minimum the data contains the following variables: a variable uniquely
76 identifying each tree-stem, `sp` of type `fct` factor identifying species, `dbh1` and `dbh2` of
77 type `dbl` quantifying the dbh at earlier and later census, and `growth` of type `dbl` double
78 quantifying the average annual growth in centimeters.

79 We load both 2008 and 2014 SCBI census data `.csv` files as they existed on GitHub on
80 November 20, 2020. After selecting only relevant variables, we perform a few additional

81 data wrangling steps: convert the variable with the date of measurement to be of type `date`,
82 convert dbh to be in centimeters¹, convert the `sp` variable containing species information
83 from type `chr` character to `fct` factor (we will discuss the need for this in Section 2.5).
84 Furthermore, in order to speed up computation for purposes of this example, we only
85 consider a 9 ha subsection of the 25.6 ha of the SCBI site: `gx` from 0–300 instead of 0–400
86 and `gy` from 300–600 instead of 0–640.

```
census_2013_scbi <- read_csv("scbi.stem2.csv") %>%  
  select(stemID, sp, date = ExactDate, gx, gy, dbh, codes, status) %>%  
  mutate(  
    date = mdy(date),  
    dbh = as.numeric(dbh)/10,  
    sp = factor(sp)  
  ) %>%  
  filter(gx < 300, between(gy, 300, 600))  
  
census_2018_scbi <- read_csv("scbi.stem3.csv") %>%  
  select(stemID, sp, date = ExactDate, gx, gy, dbh, codes, status) %>%  
  mutate(  
    date = mdy(date),  
    dbh = as.numeric(dbh)/10,  
    sp = factor(sp)  
  ) %>%  
  filter(gx < 300, between(gy, 300, 600))
```

87 These two data frames are then used as the two primary arguments to the `compute_growth()`
88 function, along with the `id` argument whereby the user specifies the name of the variable
89 that uniquely identifies each tree-stem under consideration (note this does not include
90 resprouts in the later census):

¹A rule of thumb to determine the units of dbh is check if the smallest non-zero and non-missing measurement is 1 or 10. If the former, then centimeters. If the latter, then millimeters. This is because ForestGEO protocols state that only trees with dbh greater or equal to 1cm should be included in censuses.

```

growth_scbi <-
  compute_growth(
    census_1 = census_2013_scbi,
    census_2 = census_2018_scbi %>% filter(!str_detect(codes, "R")),
    id = "stemID"
  )

growth_scbi
## Simple feature collection with 7954 features and 8 fields
## geometry type: POINT
## dimension: XY
## bbox: xmin: 0.2 ymin: 300 xmax: 299.9 ymax: 600
## CRS: NA
## # A tibble: 7,954 x 9
##   stemID sp      dbh1 codes1 status  dbh2 codes2 growth   geometry
##   <dbl> <fct> <dbl> <dbl> <chr>  <dbl> <dbl> <dbl> <POINT>
## 1 4 nysy  13.6   M     A     14.2   M    0.103 (14.2 428.5)
## 2 5 havi  8.8    M     A     9.6    M;P  0.150 (9.4 436.4)
## 3 6 havi  3.25   NULL  A     4      M    0.140 (1.3 434)
## 4 77 qual 65.2   M     A     66     M    0.141 (34.7 307.2)
## 5 79 tiam  47.7   M     A     46.8   M   -0.161 (40 381.1)
## 6 80 caca  5.15   M     A     6.5    M    0.253 (38.7 421.7)
## 7 96 libe  2.3    J;M   A     3.7    M    0.262 (60 310)
## 8 100 caca 5.09   NULL  A     NA     DN   NA     (52.5 476.3)
## 9 101 litu  65.4   M     A     68.4   M    0.552 (47.1 567.3)
## 10 102 astr  1.99   NULL  A     2.5    M    0.0954 (40.8 575.5)
## # ... with 7,944 more rows

```

91 The output `growth_scbi` is a single data frame of class `sf` that includes a numerical
 92 `growth` reflecting the average annual growth in dbh (in cm) for all trees that were alive at
 93 both time points as well a `geometry` variable encoding each tree's geolocation. Furthermore,
 94 variables that (in theory) remain unchanged between censuses appear only once, such as

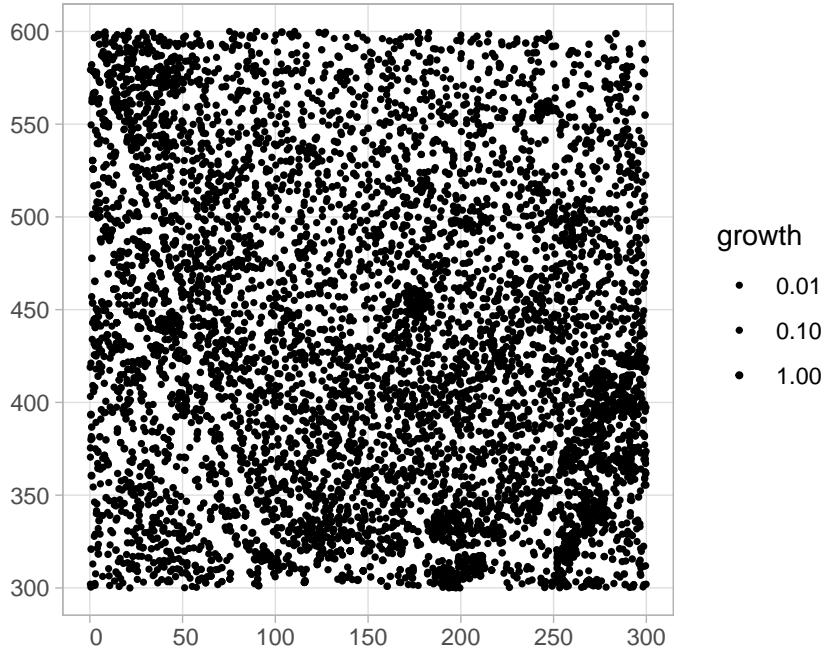


Figure 1: Growth of trees at SCBI.

95 location variables `gx` and `gy`; as well as species-related variables. Variables that should
 96 change between censuses are suffixed with 1 and 2 indicating the earlier and later censuses,
 97 such as `dbh1/dbh2` and `codes1/codes2`.

98 Given that `growth_scbi` is of class `sf`, it can be easily plotted in `ggplot2` using the
 99 `geom_sf()` geometry as seen in Figure 1.

100 TODO: Rescale points in this plot:

```
ggplot() +
  geom_sf(data = growth_scbi, aes(size = growth)) +
  scale_size(breaks = c(0.01, 0.1, 1), range = c(0.1, 1))
```

101 2.2 Add spatial information

102 We now encode spatial information to the `growth_df` data frames. First, in order to control
 103 for study region edge effects, we add “buffers” to the periphery of the study region (cite
 104 Waller?). Our model of interspecific competition relies on a spatial definition of who the
 105 competitor trees are for focal trees of interest. Since certain explanatory variables such as

106 basal area are cumulative, we must ensure that all trees being modeled are not biased to
107 have different neighbor structures. This is a particular concern for trees at the boundary
108 of study regions, which will not have the same number of neighbors as trees in the internal
109 part of the study region.

110 Second, our ultimate method for model assessment will rely on estimates of model error
111 as generated by cross-validation. Conventional cross-validation schemes assign observations
112 to folds by resampling individual observations at random. However, underlying this scheme
113 is an assumption that the observations are independent. In the case of forest census data,
114 observations exhibit spatial autocorrelation, and thus this dependence must be incorporated
115 in our resampling scheme in spatial cross-validation Roberts et al. (2017) Pohjankukka et al.
116 (2017) We will therefore associate portions of the study region to spatial folds.

117 To these two ends, we define two constants, both of which are in the same units as the
118 `gx` and `gy` variables (most often meters).

```
comp_dist <- 7.5
cv_fold_size <- 100
```

119 The first constant is `comp_dist` which defines the maximum distance for a tree's compet-
120 itive neighborhood. Trees within this distance of each other are assumed to compete while
121 those farther than this distance apart do not. Put differently, all trees within `comp_dist` of
122 a focal tree will be considered its competitors (see below). Other studies have estimated the
123 value of `comp_dist`; we use an average of estimated values Canham et al. (2004), Uriarte
124 et al. (2004), Tatsumi et al. (2013), Canham et al. (2006).

125 Furthermore, `comp_dist` will define the size of all buffers considered, which will be
126 encoded as a binary variable `buffer` as computed by the `add_buffer_variable()` function.
127 This function takes as input the main `growth_df` data frame, the `size` of the buffer which
128 we set as `comp_dist`, and the boundary of the study region encoded as a simple features
129 polygon Pebesma (2018). DESCRIBE SF PACKAGE. In the Big Woods example below
130 we will use a pre-loaded simple features polygon while for the SCBI example we present
131 example code on how to manually construct one.

132 The second constant is `cv_fold_size` which defines the length and width of the spatial
133 folds (note that for now the spatial folds are restricted be squares). We will then use this

134 constant to associate each observed tree to one of k folds in the respective study region. In
135 the Big Woods example below we will use the `blockCV` R package that has implemented
136 spatial cross-validation while for the SCBI we will do this manually Valavi et al. (2019)

137 **2.2.1 Big Woods**

138 First, we indicate which trees are part of the buffer. This necessitates information about
139 the study region boundary. In this case, we use a `sf_polygon` object `study_region_bw`
140 which comes pre-loaded in the `forestecology` packages. After loading `study_region_bw`,
141 we illustrate the results of the `add_buffer_variable()` function in Figure 2. Trees on the
142 periphery denote with lighter colors are part of the buffer and will not be considered as
143 “focal” trees of interest going forward; they will only be considered as competitor trees.

```
data(study_region_bw)

growth_bw <- growth_bw %>%
  add_buffer_variable(direction = "in", size = comp_dist, region = study_region_bw)

ggplot() +
  geom_sf(data = growth_bw %>% sample_frac(0.2), aes(col = buffer), size = 0.5)
```

144 Second, we associate each tree to spatial cross validation folds. In this case, we use the
145 `spatialBlock()` function from the `blockCV` package to define the spatial grid which
146 THIS IS A MESS. We use the Valavi et al. (2019), whose elements will act as the folds
147 in our leave-one-out (by “one” we mean “one grid block”) cross-validation scheme. The
148 upshot here is we add `foldID` to `growth_df` which identifies which fold each individual is
149 in, and the creation of a `cv_grid_sf` object which gives the geometry of the cross validation
150 grid.

```
set.seed(76)
bw_spatialBlock <- spatialBlock(
  speciesData = growth_bw, theRange = cv_fold_size, k = 28, x0ffset = 0.5,
```

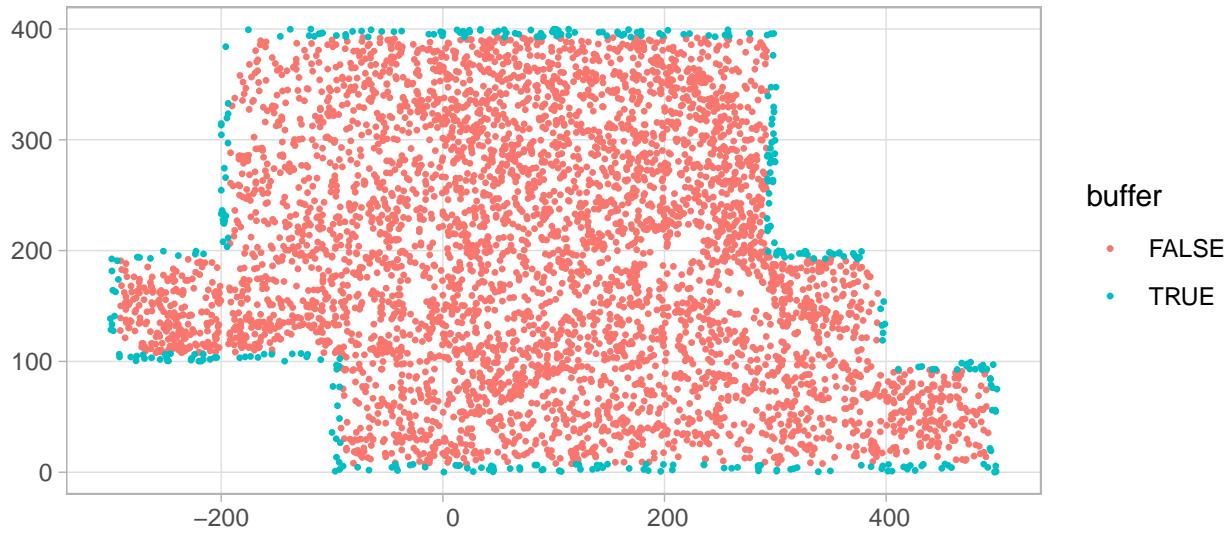


Figure 2: Buffer region for Big Woods study region.

```
yOffset = 0, verbose = FALSE, showBlocks = FALSE
)
```

151 Then add foldID to each tree

```
growth_bw <- growth_bw %>%
  mutate(foldID = bw_spatialBlock$foldID)
```

```
# Visualize grid. Why does fold 19 repeat?
ggplot() +
  geom_sf(data = bw_spatialBlock$blocks %>% st_as_sf()) +
  geom_sf(data = growth_bw %>% sample_frac(0.2),
         aes(col = factor(foldID)), size = 0.1, show.legend = FALSE) +
  geom_sf_text(data = bw_spatialBlock$blocks %>% st_as_sf(),
               aes(label = folds))
```

152 Then remove empty folds

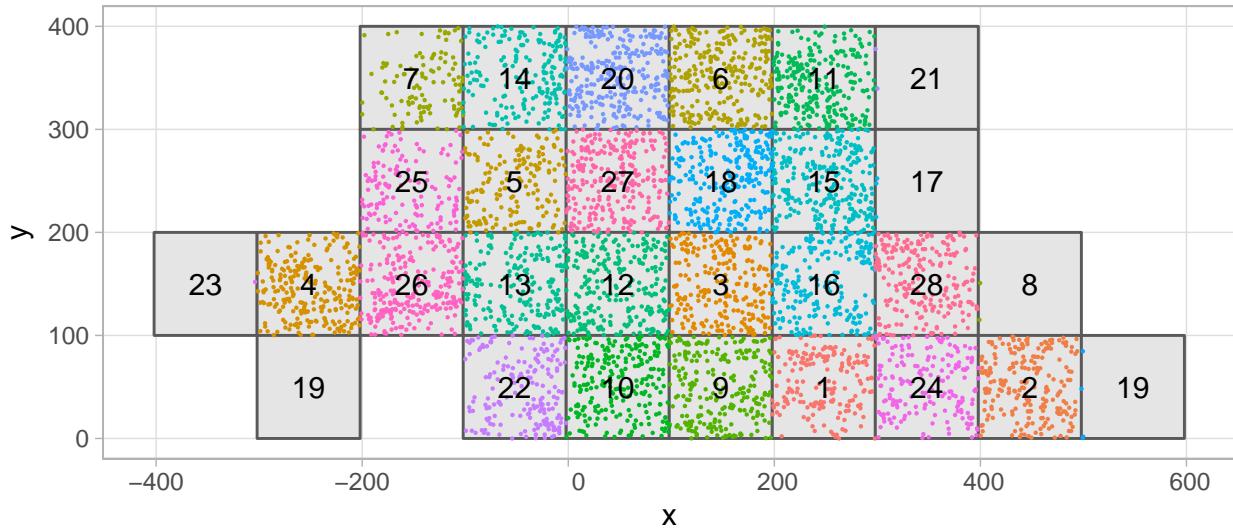


Figure 3: Inspect blocks closely.

```

growth_bw <- growth_bw %>%
  filter(!foldID %in% c(19, 23, 21, 17, 8, 19)) %>%
  mutate(foldID = factor(foldID))

```

153 Separately, we save the spatial cross-validation grid as an `sf_polygon` object `blocks_bw`

```

blocks_bw <- bw_spatialBlock$blocks %>%
  st_as_sf()

```

154 **2.2.2 SCBI**

155 First, we indicate which trees are part of the buffer. In this case however we manually define
 156 the study region boundary based on the subregion we defined in Section ?? and create an
 157 `sf_polygon` object using the `sf_polygon()` function from the `sfheaders` package. Figure
 158 4 displays the resulting buffer trees.

```

study_region_scbi <- tibble(
  x = c(0, 300, 300, 0, 0),
  y = c(300, 300, 600, 600, 300)
) %>%

```

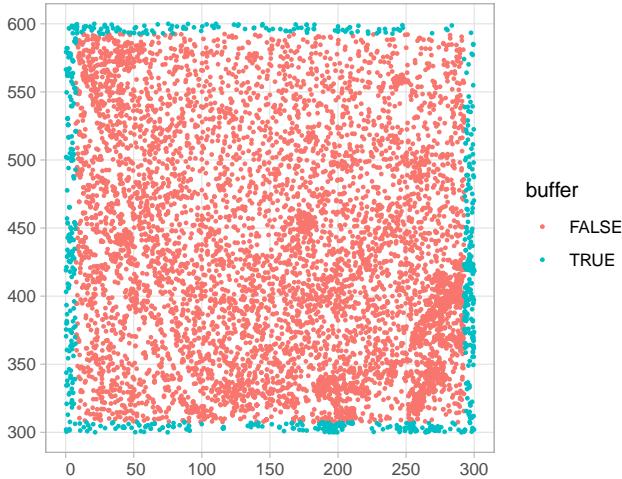


Figure 4: Buffer region for SCBI study region.

```

sf_polygon()

growth_scbi <- growth_scbi %>%
  add_buffer_variable(direction = "in", size = comp_dist, region = study_region_scbi)

ggplot() +
  geom_sf(data = growth_scbi, aes(col = buffer), size = 0.5)

```

159 Second, we associate each tree to spatial cross validation folds. In this case we manually
 160 define a spatial crossvalidation grid. Figure 5 displays the resulting cross-validation folds
 161 along with the buffer from Figure 4.

162 Here we manually define the spatial cross-validation grid as an `sf_polygon` object
 163 `scbi_cv_grid`

```

fold1 <- rbind(c(0, 300), c(150, 300), c(150, 600), c(0, 600))
fold2 <- rbind(c(150, 300), c(300, 300), c(300, 600), c(150, 600))

blocks_scbi <- bind_rows(
  sf_polygon(fold1),
  sf_polygon(fold2)

```

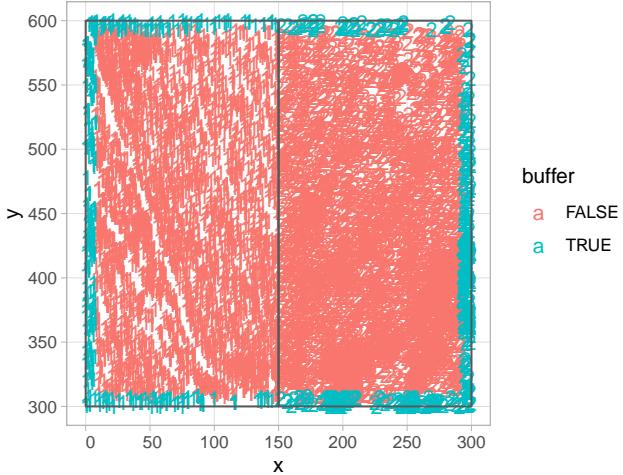


Figure 5: Buffer region for SCBI study region.

```

) %>%
  mutate(folds = c(1, 2) %>% factor())

SpatialBlock_scbi <- spatialBlock(
  speciesData = growth_scbi, k = 2, selection = "systematic", blocks = blocks_scbi,
  showBlocks = FALSE, verbose = FALSE
)

# Add foldID to each tree
growth_scbi <- growth_scbi %>%
  mutate(foldID = SpatialBlock_scbi$foldID %>% factor())

ggplot() +
  geom_sf_text(data = growth_scbi, aes(label = foldID, col = buffer)) +
  geom_sf(data = blocks_scbi, fill = "transparent")

```

164 2.3 Define focal versus competitor trees

165 Next we define `focal_vs_comp` data frames which connects each focal tree in the `growth_df`
166 data frames to the trees in its competitive neighborhood range as defined by the `comp_dist`

constant. So for example, if `growth_df` consisted of two focal trees with two and three neighbors with `comp_dist` respectively, `focal_vs_comp` would be a data frame of 5 rows connecting each focal tree to its competitors. The `create_focal_vs_comp()` function makes this connection taking as inputs the `growth_df` data frame; the `comp_dist` constant defining competitive range; `cv_grid_sf`, giving the cross validation grid; and the `id` variable.

2.3.1 Big Woods

```
focal_vs_comp_bw <- growth_bw %>%
  create_focal_vs_comp(comp_dist, cv_grid_sf = blocks_bw, id = "treeID")
```

TODO: Figure out how to show this data frame's contents.

2.3.2 SCBI

```
focal_vs_comp_scbi <- growth_scbi %>%
  create_focal_vs_comp(comp_dist, cv_grid_sf = blocks_scbi, id = "stemID")
```

TODO: Figure out how to show this data frame's contents.

2.4 Fit model and make predictions

Next we fit the following linear model to the dbh of each focal tree. Let $i = 1, \dots, n_j$ index all n_j trees of “focal” species group j ; let $j = 1, \dots, J$ index all J focal species groups; and let $k = 1, \dots, K$ index all K “competitor” species groups. We modeled the growth in diameter per year y_{ij} (in centimeters per year) of the i^{th} tree of focal species group j as a linear model f of the following covariates \vec{x}_{ij}

$$y_{ij} = f(\vec{x}_{ij}) + \epsilon_{ij} = \beta_{0,j} + \beta_{\text{DBH},j} \cdot \text{DBH}_{ij} + \sum_{k=1}^K \lambda_{jk} \cdot \text{BA}_{ijk} + \epsilon_{ij}$$

We estimate the model's parameters using Bayesian linear regression implemented in the `fit_bayesian_model()` function. TODO: define all parameters

184 For this linear model's case, there exists a closed form solution as described here. As
185 such, the `fit_bayesian_model()` function using matrix algebra to obtain all parameter
186 estimates, rather than computationally expensive Monte Carlo approximations. The inputs
187 to this function are a `focal_vs_comp` data frame, `prior_param` a list of priors, and a boolean
188 flag `run_shuffle` on whether or not to run competitor-species identity permutations which
189 we will demonstrate below on the Michigan Big Woods data. This function returns the
190 posterior means of all parameters.

191 Using these posterior means, we then use the posterior predictive distribution to obtain
192 fitted/predicted values \hat{y} of the dbh for each focal tree using the `predict_bayesian_model()`.
193 These \hat{y} can then be compared to the observed y dbh's to compute the root mean-square
194 error, a measure of a model's predictive error which has the same units as the observed
195 data y .

196 2.4.1 Big Woods

197 For the Michigan Big Woods data we present two use cases of the model fitting and pre-
198 diction scheme. The first use case is the simplest where we assess the fit of the model using
199 root mean squared error. The second use case then answers the question of whether species
200 competitor identity matters using permutation test.

201 For the first use case, we fit the linear model specified in Equation XXX to our data
202 frame of type `focal_vs_comp`. This input/outputs of the `fit_bayesian_model()` function
203 are lists of the prior/posterior means of parameters of the linear regression specified in
204 XXX. Generally speaking, there are two classes of regression parameters: β main effects
205 and λ competitive effects. In the upcoming Section 2.6, we will present code visualizing
206 this posterior distributions.

```
comp_bayes_lm_bw <- focal_vs_comp_bw %>%  
  comp_bayes_lm(prior_param = NULL)
```

207 This output of posterior parameters for the specified competition model are then used
208 along with the posterior predictive distribution encoded in `predict_bayesian_model()` to
209 return predicted growths for each individual tree. We join these predicted growths to the

210 original growth data frame.

```
focal_vs_comp_bw <- focal_vs_comp_bw %>%
  mutate(growth_hat = predict(comp_bayes_lm_bw, focal_vs_comp_bw))
```

211 We then use the `rmse()` function from the `yardstick` package to obtain the root mean
212 squared error of the observed versus fitted values of growth.

```
focal_vs_comp_bw %>%
  rmse(truth = growth, estimate = growth_hat) %>%
  pull(.estimate)
## [1] 0.148145
```

213 The second use case is near identical to the first, but with a small change in the code
214 to test whether the identity of the competitor matters. By adding a `run_shuffle = TRUE`
215 argument to `fit_bayesian_model()`, for each focal tree its competitor trees' species identity
216 will be “shuffled” randomly much like in a permutation test. By shuffling these species
217 labels we are effectively fitting the model under a null model that competitor species identity
218 does not matter. If the “shuffled” RMSE’s are consistently lower than the unshuffled RMSE
219 corresponding to the observed data, then we have evidence to suggest that competitor
220 identity matters to competitive interactions.

```
comp_bayes_lm_bw_shuffle <- focal_vs_comp_bw %>%
  comp_bayes_lm(prior_param = NULL, run_shuffle = TRUE)
```

```
focal_vs_comp_bw <- focal_vs_comp_bw %>%
  mutate(growth_hat_shuffle = predict(comp_bayes_lm_bw_shuffle, focal_vs_comp_bw))

focal_vs_comp_bw %>%
  rmse(truth = growth, estimate = growth_hat_shuffle) %>%
  pull(.estimate)
## [1] 0.1505383
```

221 The RMSE is fact lower for the non-shuffled version, indicative of a better model fit.
222 This gives support for the idea that competitor identity does matter for competitive inter-
223 actions. In Allen & Kim (2020) we run this shuffle a large number of times to construct a
224 full permutation distribution to show that this difference is robust to resampling variation.

225 **2.4.2 SCBI**

226 In the case of the SCBI data, we once again perform the same model fitting and computing
227 of fitted growths as with the Big Woods data, but this time we map the residuals of the
228 observed minus fitted values to look for spatial patterns.

```
comp_bayes_lm_scbi <- focal_vs_comp_scbi %>%  
  comp_bayes_lm(prior_param = NULL)  
  
focal_vs_comp_scbi <- focal_vs_comp_scbi %>%  
  mutate(growth_hat = predict(comp_bayes_lm_scbi, focal_vs_comp_scbi))  
  
focal_vs_comp_scbi %>%  
  rmse(truth = growth, estimate = growth_hat) %>%  
  pull(.estimate)  
## [1] 0.1281398
```

229 In Figures 6 and 7 we present the residuals.

230 **2.5 Run spatial cross-validation**

231 The model fits and predictions in Section 2.4 all suffer from a common failing: they use
232 the same data to both fit the model and to assess the model's performance using the
233 RMSE. As argued by Roberts et al. (2017), this can lead to overly optimistic assessments
234 of model quality as the models can be overfit, in particular in situations where spatial-
235 autocorrelation is present. To mitigate the effects of such overfitting, we use a spatially
236 block cross-validation algorithm implemented in the `run_cv()`. This function at its core
237 uses the same model fitting implemented in the `fit_bayesian_model()` function, however

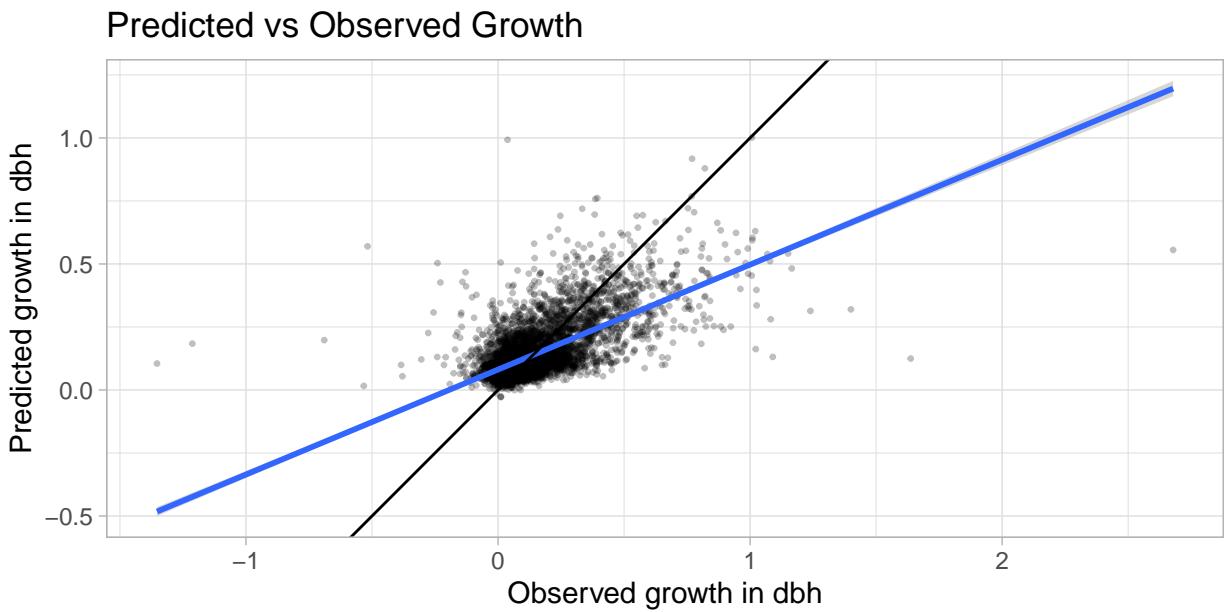


Figure 6: Predicted versus observed growth.

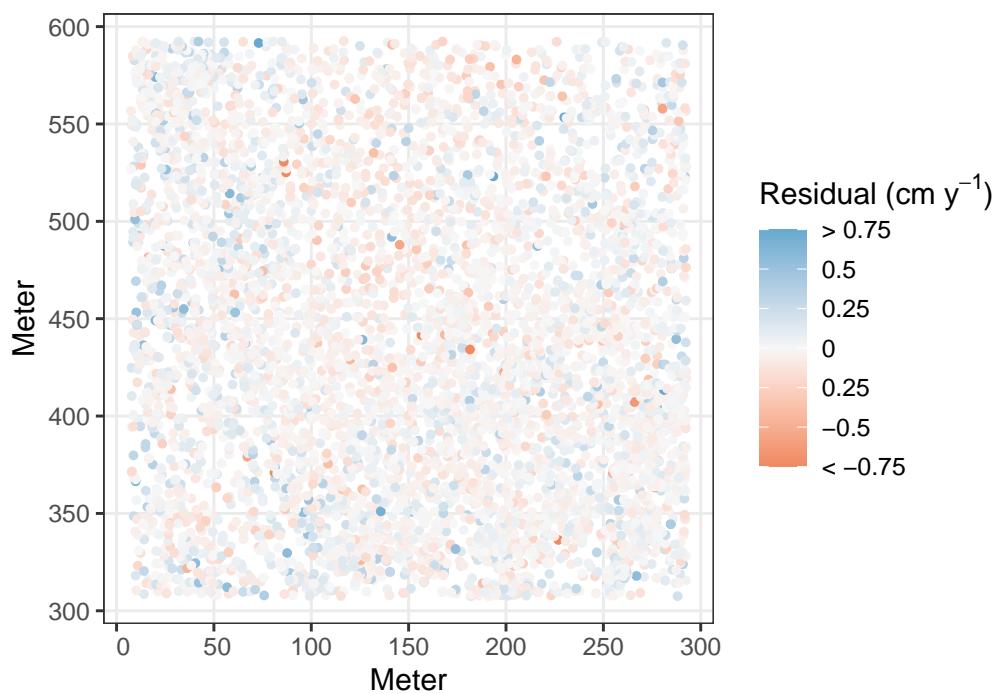


Figure 7: Spatial distribution of residuals for model applied to SCBI data.

238 trains the model on $k - 1$ spatial folds of the train and returns fitted values for the test
239 data. Recall that the spatial blocking scheme was encoded in Section 2.2.

240 **2.5.1 Big Woods**

241 Applying this spatially cross-validated model fit yields an RMSE is higher than that when
242 the model is fit without cross validation. In other words, our model fits in 2.4 were overly
243 optimistic in the model's fitting power, whereas a cross-validated results yield an estimate
244 that is closer to the truth. See Allen & Kim (2020) for more discussion of this.

```
focal_vs_comp_bw <- focal_vs_comp_bw %>%  
  run_cv(comp_dist = comp_dist, cv_grid = blocks_bw)  
  
focal_vs_comp_bw %>%  
  rmse(truth = growth, estimate = growth_hat) %>%  
  pull(.estimate)  
## [1] 0.1532316
```

245 **2.5.2 SCBI**

246 Observe once again that this RMSE is much higher than that for the above SCBI model
247 fit without cross-validation.

```
focal_vs_comp_scbi <- focal_vs_comp_scbi %>%  
  run_cv(comp_dist = comp_dist, cv_grid = blocks_scbi)  
  
focal_vs_comp_scbi %>%  
  rmse(truth = growth, estimate = growth_hat) %>%  
  pull(.estimate)  
## [1] 0.144608
```

248 2.6 Visualize posterior distributions

249 Lastly, we return to the model fits from Section 2.4 and present tools to visually explore
250 the posterior distributions of all parameters in our model. There are two main groups of
251 parameters to consider. The β coefficients tell us about how fast each species grows and
252 how this depends on DBH while the full matrix of λ values describe the competitive effects
253 between pairs of species. There is a rich literature on this matrix (cite).

254 DO WE NEED TO DESCRIBE MECHANICS? Because of the structure of the `bw_fit_model`
255 object we cannot simply draw these curves based on the posterior distribution. `bw_fit_model()`
256 gives the parameters *compared* to a baseline. This is not of direct interest. So to display
257 these parameters, as we care about them, we have to sample from the baseline distribution
258 and from the comparison one to get the posterior distribution of interest.

259 2.6.1 Big Woods

260 Here we re-run the model fit to the Big Woods data from Section 2.4, but this time use “fam-
261 ily” as the group for comparison which has. This makes the posterior distributions easier to
262 follow. Also, surprisingly, grouping by family performed just as well as grouping by species
263 Allen & Kim (2020). First we re-run `create_focal_vs_comp()` and `fit_bayesian_model()`
264 with no permutation shuffling with the grouping variable as family.

```
focal_vs_comp_bw <- growth_bw %>%  
  mutate(sp = family %>% factor()) %>%  
  create_focal_vs_comp(comp_dist = comp_dist, cv_grid_sf = blocks_bw, id = "treeID")  
  
comp_bayes_lm_bw <- focal_vs_comp_bw %>%  
  comp_bayes_lm(prior_param = NULL)
```

265 Now the posterior parameter outputs of `fit_bayesian_model()` are passed to `plot_bayesian_model_pa`
266 to generate visualizations of the posterior parameters. These visualizations are displayed
267 in Figure 5 of Allen & Kim (2020). For simplicity we only plot a subset of the species
268 families.

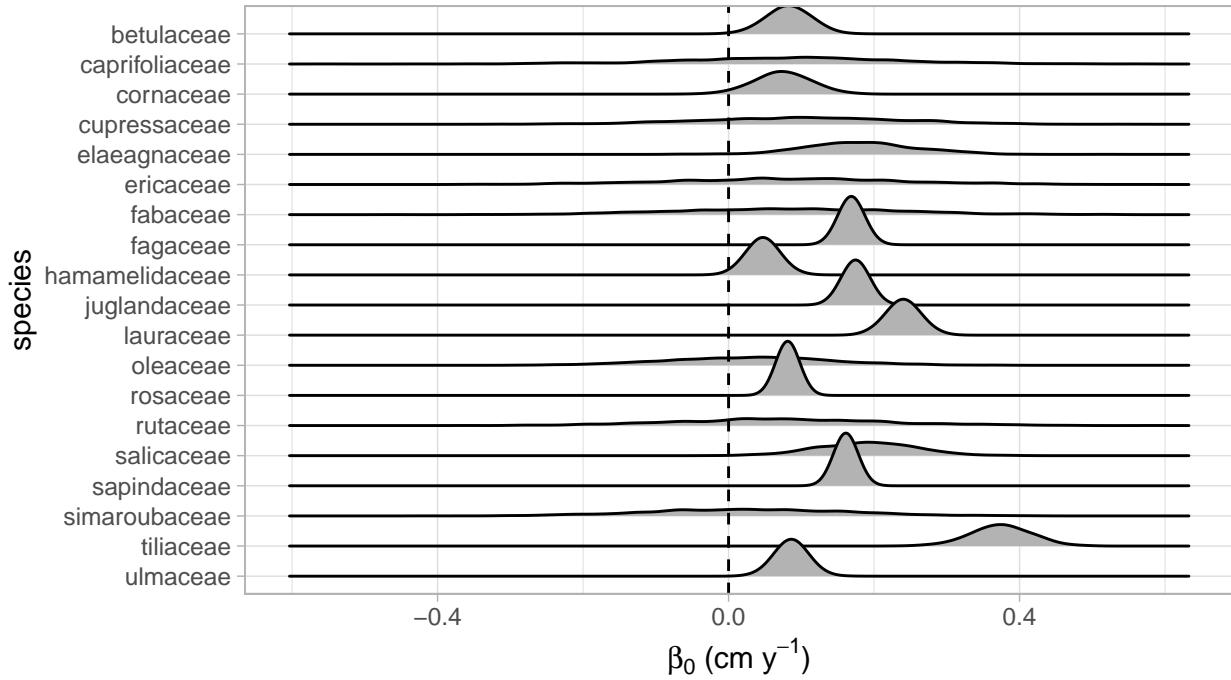


Figure 8: Posterior distribution of beta0.

```
sp_to_plot <- c("cornaceae", "fagaceae", "hamamelidaceae", "juglandaceae",
                 "lauraceae", "rosaceae", "sapindaceae", "ulmaceae")
```

269 The output is a list with three plots stored. Figure 8 The element **beta_0** gives the
 270 baseline growth intercept β_0 , i.e., how fast an individual of each group grows independent
 271 of DBH).

```
plot1 <- autoplot(comp_bayes_lm_bw, type = "intercepts")
plot1
```

272 Figure 9 Next **beta_dbh** gives the slope for DBH slope $\beta_{dbh,i}$ for each group.

```
plot2 <- autoplot(comp_bayes_lm_bw, type = "dbh_slopes")
plot2
```

273 Finally Figure 10 **lambda** gives the competition coefficients λ .

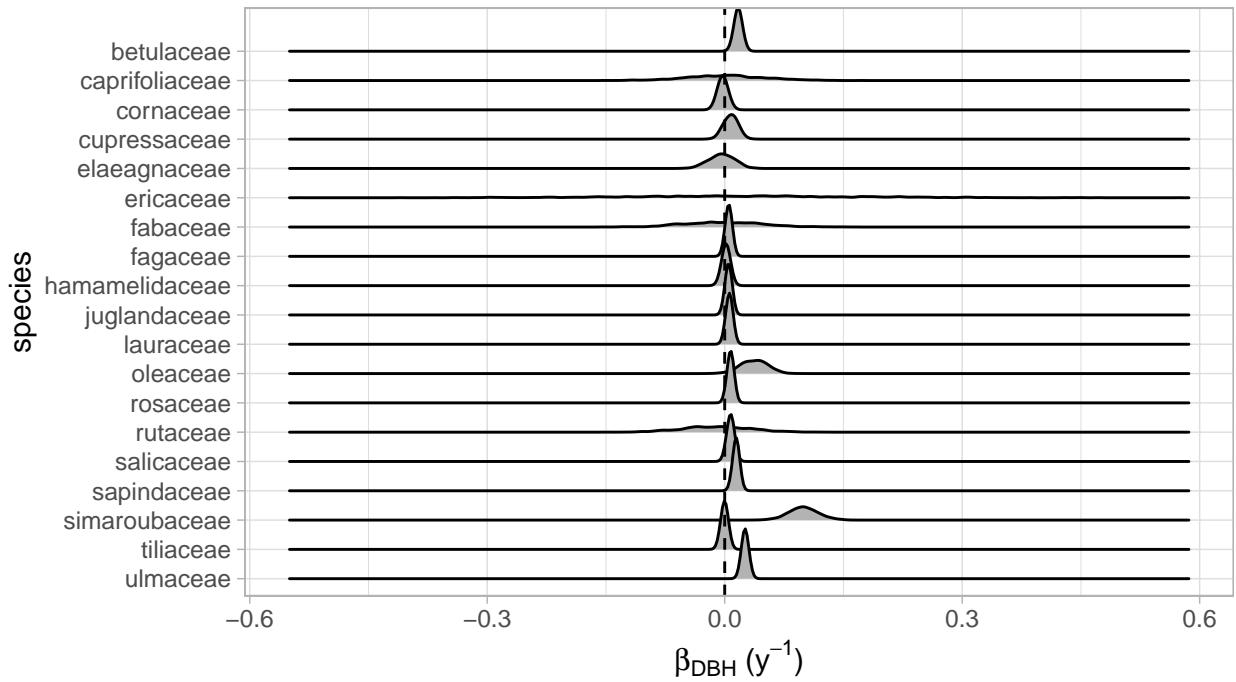


Figure 9: Posterior distribution of betadbh.

```
plot3 <- autoplot(comp_bayes_lm_bw, type = "competition")
plot3
```

274 2.6.2 SCBI

275 We revisit the posterior parameters for the SCBI from Section {model-fit-predict}, but this
 276 time only focus on the λ competition coefficients.

```
sp_to_plot <- c("quru", "litu", "cagl", "cato")
```

```
plot3 <- autoplot(comp_bayes_lm_bw, type = "competition")
plot3
```

277 Add explanation here.

278 HEY BERT PICK IT UP HERE

Competitor species in rows, focal species in columns

Ex: Top row, second column: competitive effect of betulaceae on caprifoliaceae

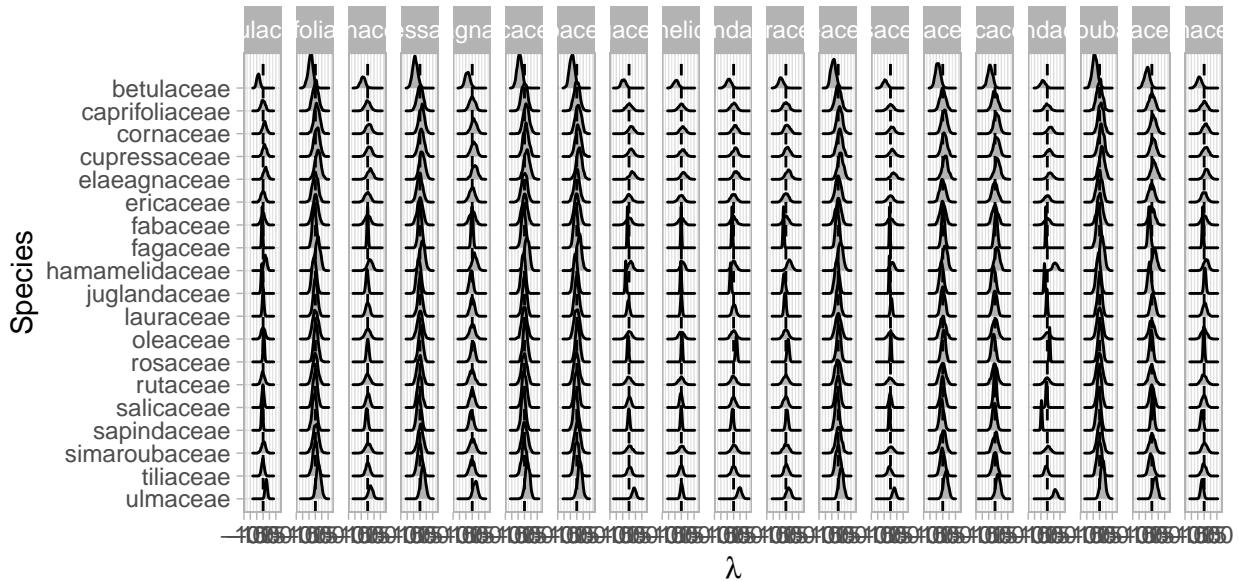


Figure 10: Posterior distribution of lambda's for Big Woods.

Competitor species in rows, focal species in columns

Ex: Top row, second column: competitive effect of betulaceae on caprifoliaceae

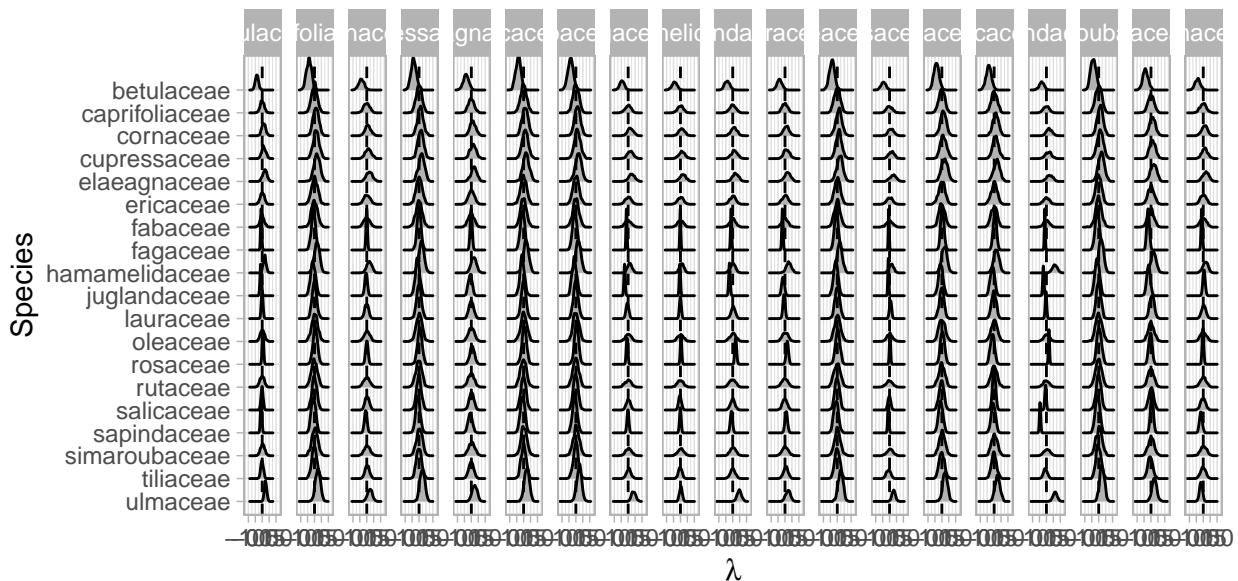


Figure 11: Posterior distribution of lambda's for SCBI.

²⁷⁹ **3 Discussion**

²⁸⁰ **4 Acknowledgments**

²⁸¹ **References**

- ²⁸² Allen, D., Dick, C., Burnham, R. J., Perfecto, I. & Vandermeer, J. (2020), 'The michigan big
²⁸³ woods research plot at the edwin s. george, pinckney, mi, usa', *Miscellaneous Publications
284 of the Museum of Zoology, University of Michigan* **207**.
- ²⁸⁵ **URL:** <http://hdl.handle.net/2027.42/156251>
- ²⁸⁶ Allen, D. & Kim, A. Y. (2020), 'A permutation test and spatial cross-validation approach
²⁸⁷ to assess models of interspecific competition between trees', *PLOS ONE* **15**(3), e0229930.
²⁸⁸ Publisher: Public Library of Science.
- ²⁸⁹ **URL:** <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0229930>
- ²⁹⁰ Anderson-Teixeira, K. J., Davies, S. J., Bennett, A. C., Gonzalez-Akre, E. B., Muller-
²⁹¹ Landau, H. C., Wright, S. J., Salim, K. A., Zambrano, A. M. A., Alonso, A., Baltzer,
²⁹² J. L., Basset, Y., Bourg, N. A., Broadbent, E. N., Brockelman, W. Y., Bunyavejchewin,
²⁹³ S., Burslem, D. F. R. P., Butt, N., Cao, M., Cardenas, D., Chuyong, G. B., Clay, K.,
²⁹⁴ Cordell, S., Dattaraja, H. S., Deng, X., Detto, M., Du, X., Duque, A., Erikson, D. L.,
²⁹⁵ Ewango, C. E. N., Fischer, G. A., Fletcher, C., Foster, R. B., Giardina, C. P., Gilbert,
²⁹⁶ G. S., Gunatilleke, N., Gunatilleke, S., Hao, Z., Hargrove, W. W., Hart, T. B., Hau, B.
²⁹⁷ C. H., He, F., Hoffman, F. M., Howe, R. W., Hubbell, S. P., Inman-Narahari, F. M.,
²⁹⁸ Jansen, P. A., Jiang, M., Johnson, D. J., Kanzaki, M., Kassim, A. R., Kenfack, D.,
²⁹⁹ Kibet, S., Kinnaird, M. F., Korte, L., Kral, K., Kumar, J., Larson, A. J., Li, Y., Li, X.,
³⁰⁰ Liu, S., Lum, S. K. Y., Lutz, J. A., Ma, K., Maddalena, D. M., Makana, J.-R., Malhi,
³⁰¹ Y., Marthews, T., Serudin, R. M., McMahon, S. M., McShea, W. J., Memiaghe, H. R.,
³⁰² Mi, X., Mizuno, T., Morecroft, M., Myers, J. A., Novotny, V., Oliveira, A. A. d., Ong,
³⁰³ P. S., Orwig, D. A., Ostertag, R., Ouden, J. d., Parker, G. G., Phillips, R. P., Sack, L.,
³⁰⁴ Sainge, M. N., Sang, W., Sri-ngernyuang, K., Sukumar, R., Sun, I.-F., Sungpalee, W.,
³⁰⁵ Suresh, H. S., Tan, S., Thomas, S. C., Thomas, D. W., Thompson, J., Turner, B. L.,

- 306 Uriarte, M., Valencia, R., Vallejo, M. I., Vicentini, A., Vrška, T., Wang, X., Wang, X.,
307 Weiblen, G., Wolf, A., Xu, H., Yap, S. & Zimmerman, J. (2015), ‘CTFS-ForestGEO: a
308 worldwide network monitoring forests in an era of global change’, *Global Change Biology*
309 **21**(2), 528–549.
- 310 **URL:** <http://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.12712>
- 311 Bourg, N. A., McShea, W. J., Thompson, J. R., McGarvey, J. C. & Shen, X. (2013), ‘Initial
312 census, woody seedling, seed rain, and stand structure data for the SCBI SIGEO Large
313 Forest Dynamics Plot’, *Ecology* **94**(9), 2111–2112.
- 314 **URL:** <http://esajournals.onlinelibrary.wiley.com/doi/abs/10.1890/13-0010.1>
- 315 Canham, C. D., LePage, P. T. & Coates, K. D. (2004), ‘A neighborhood analysis of canopy
316 tree competition: effects of shading versus crowding’, *Canadian Journal of Forest Re-*
317 *search* **34**(4). Publisher: NRC Research Press Ottawa, Canada.
- 318 **URL:** <https://cdnsciencepub.com/doi/abs/10.1139/x03-232>
- 319 Canham, C. D., Papaik, M. J., Uriarte, M., McWilliams, W. H., Jenkins, J. C.
320 & Twery, M. J. (2006), ‘Neighborhood Analyses Of Canopy Tree Competi-
321 tion Along Environmental Gradients In New England Forests’, *Ecological Applica-*
322 *tions* **16**(2), 540–554. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1890/1051-0761%282006%29016%5B0540%3ANAOCTC%5D2.0.CO%3B2>.
- 324 Pebesma, E. (2018), ‘Simple Features for R: Standardized Support for Spatial Vector Data’,
325 *The R Journal* **10**(1), 439–446.
- 326 **URL:** <https://journal.r-project.org/archive/2018/RJ-2018-009/index.html>
- 327 Pebesma, E. J. & Bivand, R. S. (2005), ‘Classes and methods for spatial data in R’, *R*
328 *News* **5**(2), 9–13.
- 329 **URL:** <https://CRAN.R-project.org/doc/Rnews/>
- 330 Pohjankukka, J., Pahikkala, T., Nevalainen, P. & Heikkonen, J. (2017), ‘Estimating the
331 prediction performance of spatial models via spatial k-fold cross validation’, *International*
332 *Journal of Geographical Information Science* **31**(10), 2001–2019.

- 333 Roberts, D. R., Bahn, V., Ciuti, S., Boyce, M. S., Elith, J., Guillera-Arroita, G., Hauen-
334 stein, S., Lahoz-Monfort, J. J., Schröder, B., Thuiller, W., Warton, D. I., Wintle, B. A.,
335 Hartig, F. & Dormann, C. F. (2017), ‘Cross-validation strategies for data with temporal,
336 spatial, hierarchical, or phylogenetic structure’, *Ecography* **40**(8), 913–929.
337 **URL:** <http://onlinelibrary.wiley.com/doi/abs/10.1111/ecog.02881>
- 338 Smith, W. B. (2002), ‘Forest inventory and analysis: a national inventory and monitoring
339 program’, *Environmental pollution* **116**, S233–S242.
- 340 Tatsumi, S., Owari, T., Ohkawa, A. & Nakagawa, Y. (2013), ‘Bayesian modeling of neigh-
341 borhood competition in uneven-aged mixed-species stands’, *Formath* **12**, 191–209.
- 342 Uriarte, M., Condit, R., Canham, C. D. & Hubbell, S. P. (2004), ‘A spa-
343 tially explicit model of sapling growth in a tropical forest: does the iden-
344 tity of neighbours matter?’, *Journal of Ecology* **92**(2), 348–360. _eprint:
345 <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.0022-0477.2004.00867.x>.
346 **URL:** <http://besjournals.onlinelibrary.wiley.com/doi/abs/10.1111/j.0022-0477.2004.00867.x>
347
- 348 Valavi, R., Elith, J., Lahoz-Monfort, J. J. & Guillera-Arroita, G. (2019), ‘blockCV: An
349 r package for generating spatially or environmentally separated folds for k-fold cross-
350 validation of species distribution models’, *Methods in Ecology and Evolution* **10**(2), 225–
351 232. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/2041-210X.13107>.
352 **URL:** <http://besjournals.onlinelibrary.wiley.com/doi/abs/10.1111/2041-210X.13107>
- 353 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., Grole-
354 mund, G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T. L., Miller, E., Bache,
355 S. M., Müller, K., Ooms, J., Robinson, D., Seidel, D. P., Spinu, V., Takahashi, K.,
356 Vaughan, D., Wilke, C., Woo, K. & Yutani, H. (2019), ‘Welcome to the Tidyverse’,
357 *Journal of Open Source Software* **4**(43), 1686.
358 **URL:** <https://joss.theoj.org/papers/10.21105/joss.01686>