

forestecology 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 2, 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 1 Abstract (350 words)

² 1. Set the context for and purpose of the work: The scientific question/problem and the
³ desiderata

- ⁴ • (Eventually) modularly fitting models for interspecific competition and assessing
⁵ them using spatial crossvalidation
- ⁶ • Leverage ForestGEO protocols providing standardization

⁷ 2. Indicate the approach and the methods: Use tidyverse, simple features, and tidy-
⁸ models packages.

- ⁹ • tidyverse: as stated in tidy tools manifesto: standadized data structures, func-
¹⁰ tional programming with pipe, designed for humans
- ¹¹ • sf: tidyverse-friendly package that makes wrangling and visualizing spatial data
¹² much easier
- ¹³ • tidymodels: given our spatial-crossvalidation, use tidymodels framework is a col-
¹⁴ lection of packages for modeling and machine learning using tidyverse principles
- ¹⁵ • Dave: hmm, I think here we might need to mention the methods we implement.
¹⁶ Something like, the package provides functions to specific a linear, bayesian
¹⁷ neighborhood competition model of growth, then fit the model, and compare
¹⁸ competing models with spatiall cross validation. Or somethign like that.

¹⁹ 3. Outline the main results: We replicate from scratch the figure in PLOSOne paper
²⁰ using Big Woods data, conduct similar analysis for SCBI data.

- ²¹ • Code would've been way more complicated in base.
- ²² • We don't have to worry about how the component functions work, only that
²³ sequence/converyor belt is correct and output is correct.
- ²⁴ • Scientist is abstracted away from ugly programming details.

²⁵ 4. Identify the conclusions and wider implications:

- ²⁶ • New scientific conclusions from SCBI data
- ²⁷ • Modularly switch out our bayesian lm() functions to anything you want.
- ²⁸ • this can serve as blue print for other modeling situations.

²⁹ **2 Introduction**

³⁰ Repeat-censused forest plots offer excellent data to test neighborhood models of tree com-
³¹ petition Allen & Kim (2020) Canham et al. (2006) Uriarte et al. (2004). Here we describe
³² an R package, **forestecology**, to do that. This package implements the methods in Allen
³³ & Kim (2020). It provides: a convenient way to specify and fit models of tree growth based
³⁴ on neighborhood competition; a spatial cross validation method to test and compare model
³⁵ fits Roberts et al. (2017); and an ANOVA-like method to assess whether the competitor
³⁶ identity matters in these models. The model is written to work with ForestGEO plot data
³⁷ Anderson-Teixeira et al. (2015), but we envision that it could easily be modified to work
³⁸ with data from other forest plots, e.g. the US Forest Service Forest Inventory and Analysis
³⁹ plots Smith (2002).

⁴⁰ **3 Example**

⁴¹ We demonstrate the **forestecology** package's features on two data sets, both based on
⁴² inventory censuses of two sites from the Smithsonian Institution's ForestGEO international
⁴³ network of 72 long-term forest dynamics research sites Anderson-Teixeira et al. (2015).
⁴⁴ First, the Michigan Big Woods Forest Dynamics Plot located at the Edwin S. George
⁴⁵ Reserve in Pinckney, MI, USA. The 23 ha plot is situated in mature oak-hickory forest.
⁴⁶ The canopy is dominated by white oak (*Quercus alba*), northern red oak (*Q. rubra*), black
⁴⁷ oak (*Q. velutina*), shagbark hickory (*Carya ovata*) and pignut hickory (*C. glabra*) Allen
⁴⁸ et al. (2020). In the example below, we will preface any data frames from this plot in with
⁴⁹ **bw_**.

⁵⁰ Second, the Smithsonian Conservation Biology Institute (SCBI) large forest dynamics
⁵¹ plot, located at the Smithsonian's National Zoo and Conservation Biology Institute in
⁵² Front Royal, VA, USA. The 25.6 ha (640 x 400 m) plot is located at the intersection of
⁵³ three of the major physiographic provinces of the eastern US: the Blue Ridge, Ridge and
⁵⁴ Valley, and Piedmont provinces and is adjacent to the northern end of Shenandoah National
⁵⁵ Park. The forest type is typical mature secondary eastern mixed deciduous forest, with
⁵⁶ a canopy dominated by tulip poplar (*Liriodendron tulipifera*), oaks (*Quercus* spp.), and

57 hickories (*Carya* spp.), and an understory composed mainly of spicebush (*Lindera benzoin*),
58 paw-paw (*Asimina triloba*), American hornbeam (*Carpinus caroliniana*), and witch hazel
59 (*Hamamelis virginiana*) Bourg et al. (2013). In the example below, we will preface any
60 data frames from this plot in with `scbi_`.

61 The code that generates Figures are included in the supplementary materials.

62 We load all the necessary packages.

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

63 3.1 Preprocess census data

64 We start by preprocessing the census data for both sites. While ForestGEO data protocols
65 ensure a high degree of standardization between site, minor variations still exist Anderson-
66 Teixeira et al. (2015). While the Big Woods data comes pre-loaded in the `forestecology`
67 package, we load the SCBI data as they are saved in .csv files in the SCBI-ForestGEO-Data
68 repository on GitHub Gonzalez-Akre et al. (2020). In both cases, we load the census data
69 as R as “tibble” data frames thereby ensuring a standardized input/output format that
70 can be used across all `tidyverse` packages Wickham et al. (2019).

71 Furthermore, we ensure that the different variables have the correct names, types (`dbl`,
72 `data, factor`).

73 3.1.1 Big Woods

74 We load census data from 2008 and 2014 saved in the package, then merge species data
75 (genus, species, linnean classification, family, etc).

```
census_2008_bw <- census_2008_bw %>%
  left_join(species_bw, by = "sp") %>%
  select(-c(genus, species, latin))
```

76 3.1.2 SCBI

77 We load census data from 2008 and 2014 from .csv files saved from GitHub on November
78 20, 2020. Furthermore, we perform two additional pre-processing steps. First, in order to
79 speed up computation for purposes of this example, we only consider a 9 ha subsection of
80 the 25.6 ha of the SCBI site: gx from 0–300 instead of 0–400 and gy from 300–600 instead
81 of 0–640. Second, in order to standardize comparisons between Big Woods and SCBI, we
82 convert the units of dbh from mm to cm.¹

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

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

83 **3.2 Compute annual growth**

84 MERGE THIS SUBSECTION with previous?

85 For each plot we then compute average annual growth between the two censuses using
86 the `compute_growth()` function. This function takes the two census data frames as well as a
87 character indicating which variable in both data frames uniquely identifies each stem. This
88 function returns a single data frame that includes a numerical variable `growth` reflecting
89 the average annual dbh growth (in cm) of all trees alive at both time points. Furthermore,
90 variables that (in theory) remain unchanged between censuses appear only once, such as
91 location variables `gx` and `gy`; as well as species-related variables. Variables that should
92 change between censuses are suffixed with 1 and 2 indicating the earlier and later censuses,
93 such as `dbh1/dbh2` and `codes1/codes2`. Here the resulting data frames are named with
94 some variation of `growth_df`.

95 After computing the average annual growth for each tree, we ensure to convert all
96 variables denote species from type character to factors; this is to ensure that issues of rare
97 species being accounted for in both training and test sets in our upcoming cross-validation
98 step (see Section REF)

99 **3.2.1 Big Woods**

100 In the case of Big Woods data, we first remove all trees that were re-sprouts in the later
101 (2014) census. Additionally, we have included two classification of tree species: `species`
102 and `family`. To illustrate model comparison, we test competition models in which individ-
103 uals are grouped by species and by family.

```
growth_bw <-  
  compute_growth(  
    census_1 = census_2008_bw %>%  
      mutate(sp = to_any_case(sp) %>% factor()),  
    census_2 = census_2014_bw %>%  
      filter(!str_detect(codes, "R")) %>%  
      mutate(sp = to_any_case(sp) %>% factor()),
```

```
    id = "treeID"  
)  
)
```

104 **3.2.2 SCBI**

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

105 **3.2.3 Comparison**

106 DAVE: I THINK THIS FIGURE IS UNNEEDED. COULD BE REMOVED FOR SPACE
107 REASONS. Figure 1 displays histograms comparing the distribution of average annual
108 growth at both sites. Observe that average annual growth appears higher at the Big
109 Woods site.

110 **3.3 Add spatial information**

111 We now encode spatial information to the `growth_df` data frames. First, in order to control
112 for study region edge effects, we add “buffers” to the periphery of the study region (cite
113 Waller?). Our model of interspecific competition relies on a spatial definition of who the
114 competitor trees are for focal trees of interest. Since certain explanatory variables such as
115 basal area are cumulative, we must ensure that all trees being modeled are not biased to
116 have different neighbor structures. This is a particular concern for trees at the boundary
117 of study regions, which will not have the same number of neighbors as trees in the internal
118 part of the study region.

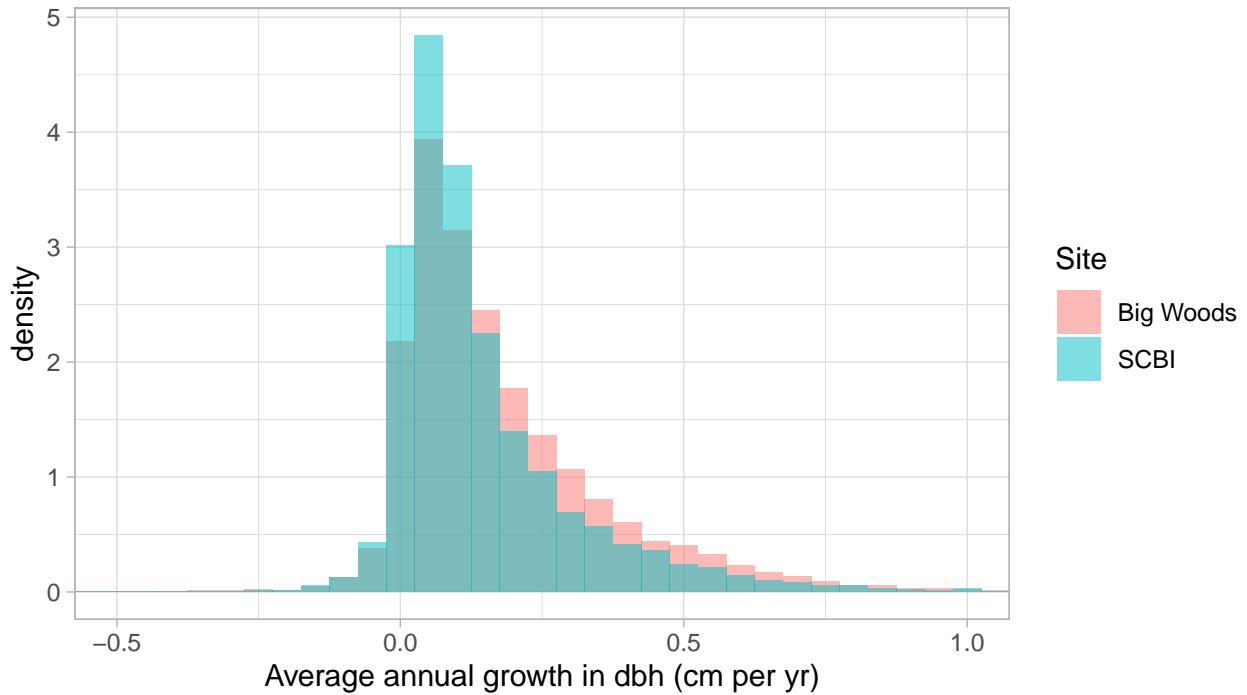


Figure 1: Distribution of average annual growth in DBH for both sites.

119 Second, our ultimate method for model assessment will rely on estimates of model error
 120 as generated by cross-validation. Conventional cross-validation schemes assign observations
 121 to folds by resampling individual observations at random. However, underlying this scheme
 122 is an assumption that the observations are independent. In the case of forest census data,
 123 observations exhibit spatial autocorrelation, and thus this dependence must be incorporated
 124 in our resampling scheme in spatial cross-validation Roberts et al. (2017) Pohjankukka et al.
 125 (2017) We will therefore associate portions of the study region to spatial folds.

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

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

128 The first constant is `comp_dist` which defines the maximum distance for a tree's compet-
 129 itive neighborhood. Trees within this distance of each other are assumed to compete while
 130 those farther than this distance apart do not. Put differently, all trees within `comp_dist` of
 131 a focal tree will be considered its competitors (see below). Other studies have estimated the

132 value of `comp_dist`; we use an average of estimated values Canham et al. (2004), Uriarte
133 et al. (2004), Tatsumi et al. (2013), Canham et al. (2006).

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

141 The second constant is `cv_fold_size` which defines the length and width of the spatial
142 folds (note that for now the spatial folds are restricted be squares). We will then use this
143 constant to associate each observed tree to one of k folds in the respective study region. In
144 the Big Woods example below we will use the `blockCV` R package that has implemented
145 spatial cross-validation while for the SCBI we will do this manually Valavi et al. (2019)

146 3.3.1 Big Woods

147 First, we indicate which trees are part of the buffer. This necessitates information about
148 the study region boundary. In this case, we use a `sf_polygon` object `study_region_bw`
149 which comes pre-loaded in the `forestecology` packages. After loading `study_region_bw`,
150 we illustrate the results of the `add_buffer_variable()` function in Figure 2. Trees on the
151 periphery denote with lighter colors are part of the buffer and will not be considered as
152 “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)
```

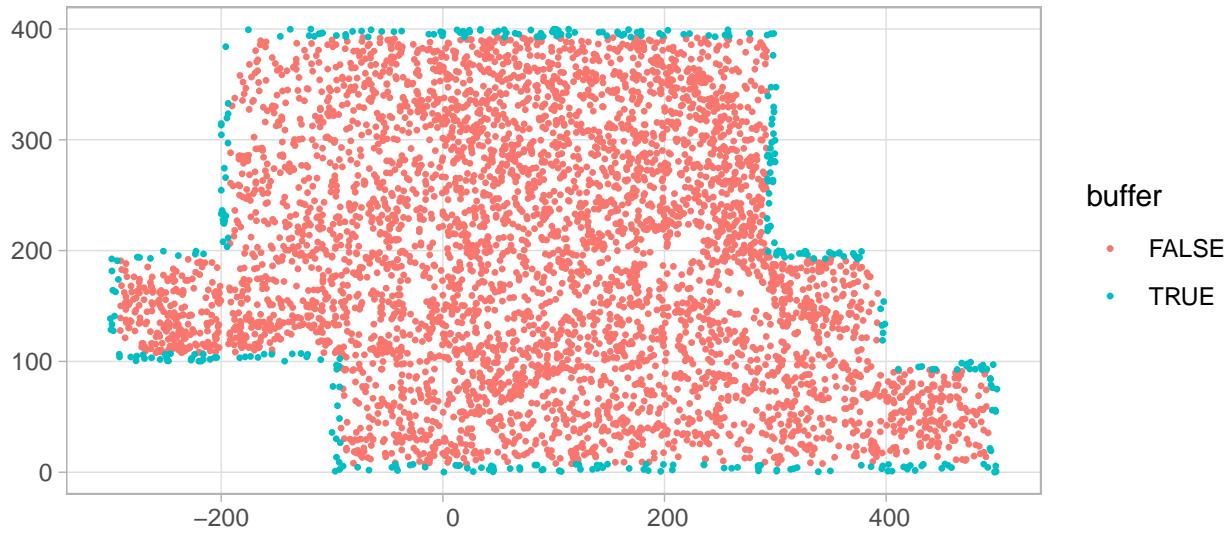


Figure 2: Buffer region for Big Woods study region.

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

```

set.seed(76)
bw_spatialBlock <- spatialBlock(
  speciesData = growth_bw, theRange = cv_fold_size, k = 28, x0ffset = 0.5,
  y0ffset = 0, verbose = FALSE, showBlocks = FALSE
)
  
```

160 Then add `foldID` to each tree

```

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

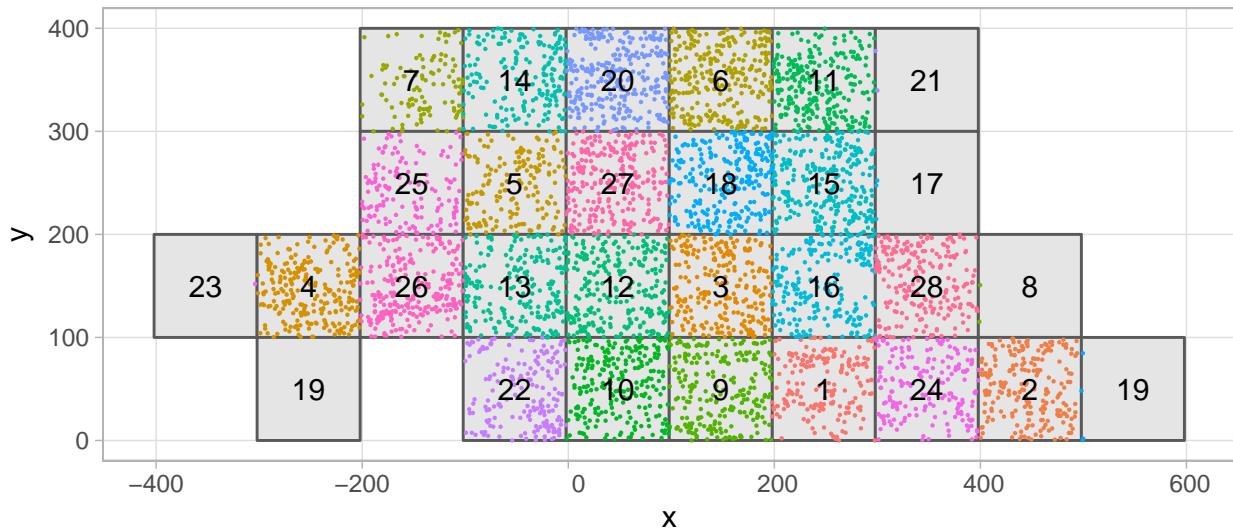


Figure 3: Inspect blocks closely.

```
# 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))
```

161 Then remove empty folds

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

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

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

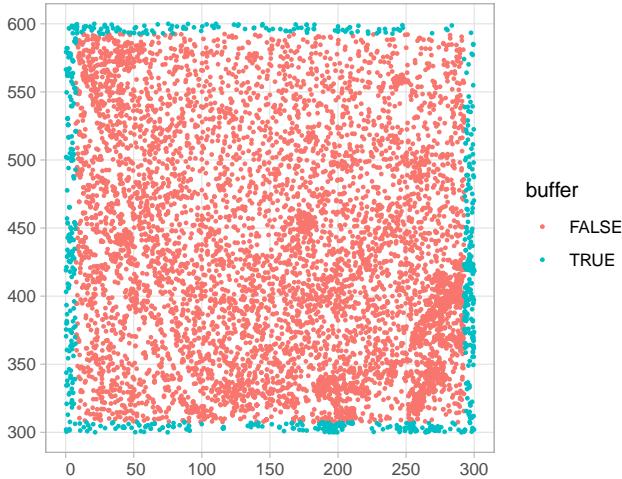


Figure 4: Buffer region for SCBI study region.

163 3.3.2 SCBI

164 First, we indicate which trees are part of the buffer. In this case however we manually define
165 the study region boundary based on the subregion we defined in Section 3.1.2 and create an
166 `sf_polygon` object using the `sf_polygon()` function from the `sfheaders` package. Figure
167 4 displays the resulting buffer trees.

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

168 Second, we associate each tree to spatial cross validation folds. In this case we manually
169 define a spatial crossvalidation grid. Figure 5 displays the resulting cross-validation folds

170 along with the buffer from Figure 4.

171 Here we manually define the spatial cross-validation grid as an `sf_polygon` object

172 `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)
) %>%
  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")
```

173 3.4 Define focal versus competitor trees

174 Next we define `focal_vs_comp` data frames which connects each focal tree in the `growth_df`
175 data frames to the trees in its competitive neighborhood range as defined by the `comp_dist`
176 constant. So for example, if `growth_df` consisted of two focal trees with two and three neigh-

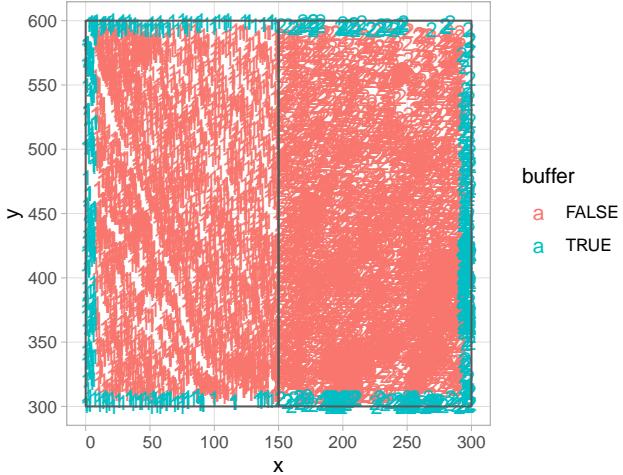


Figure 5: Buffer region for SCBI study region.

177 bors with `comp_dist` respectively, `focal_vs_comp` would be a data frame of 5 rows connect-
178 ing each focal tree to it's competitors. The `create_focal_vs_comp()` function makes this
179 connection taking as inputs the `growth_df` data frame; the `comp_dist` constant defining
180 competitive range; `cv_grid_sf`, giving the cross validation grid; and the `id` variable.

181 **3.4.1 Big Woods**

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

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

183 **3.4.2 SCBI**

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

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

185 **3.5 Fit model and make predictions**

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

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

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

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

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

205 **3.5.1 Big Woods**

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

210 For the first use case, we fit the linear model specified in Equation XXX to our data
211 frame of type `focal_vs_comp`. This input/outputs of the `fit_bayesian_model()` function

212 are lists of the prior/posterior means of parameters of the linear regression specified in
213 XXX. Generally speaking, there are two classes of regression parameters: β main effects
214 and λ competitive effects. In the upcoming Section 3.7, we will present code visualizing
215 this posterior distributions.

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

216 This output of posterior parameters for the specified competition model are then used
217 along with the posterior predictive distribution encoded in `predict_bayesian_model()` to
218 return predicted growths for each individual tree. We join these predicted growths to the
219 original growth data frame.

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

220 We then use the `rmse()` function from the `yardstick` package to obtain the root mean
221 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
```

222 The second use case is near identical to the first, but with a small change in the code
223 to test whether the identity of the competitor matters. By adding a `run_shuffle = TRUE`
224 argument to `fit_bayesian_model()`, for each focal tree its competitor trees' species identity
225 will be “shuffled” randomly much like in a permutation test. By shuffling these species
226 labels we are effectively fitting the model under a null model that competitor species identity
227 does not matter. If the “shuffled” RMSE's are consistently lower than the unshuffled RMSE
228 corresponding to the observed data, then we have evidence to suggest that competitor
229 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

```

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

234 3.5.2 SCBI

235 In the case of the SCBI data, we once again perform the same model fitting and computing
 236 of fitted growths as with the Big Woods data, but this time we map the residuals of the
 237 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

```

238 In Figures 6 and 7 we present the residuals.

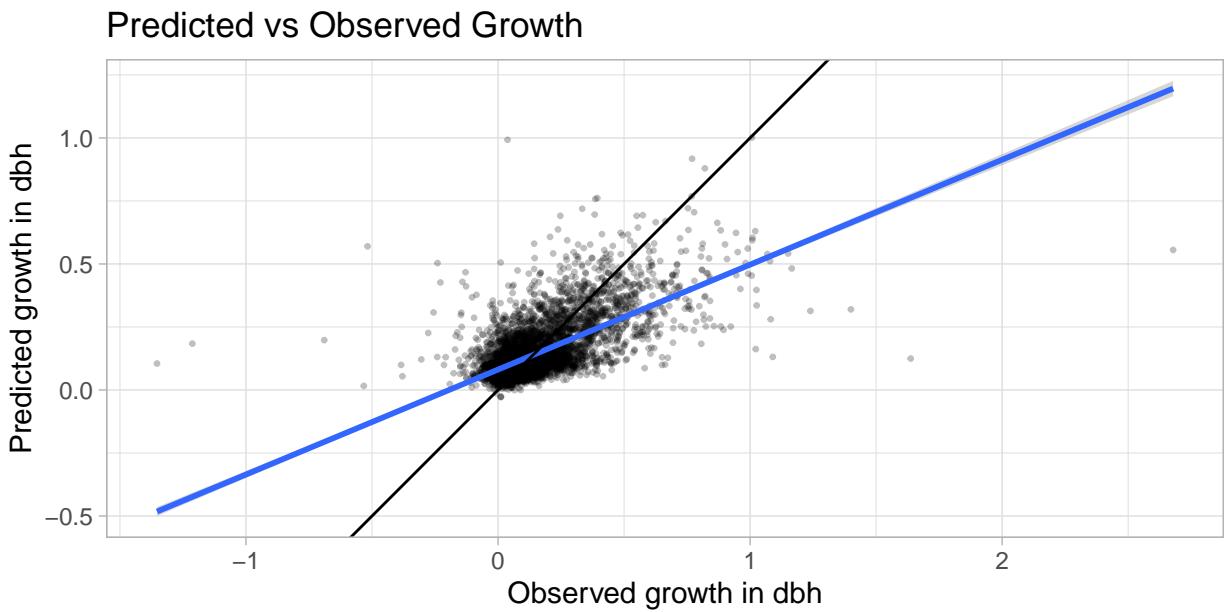


Figure 6: Predicted versus observed growth.

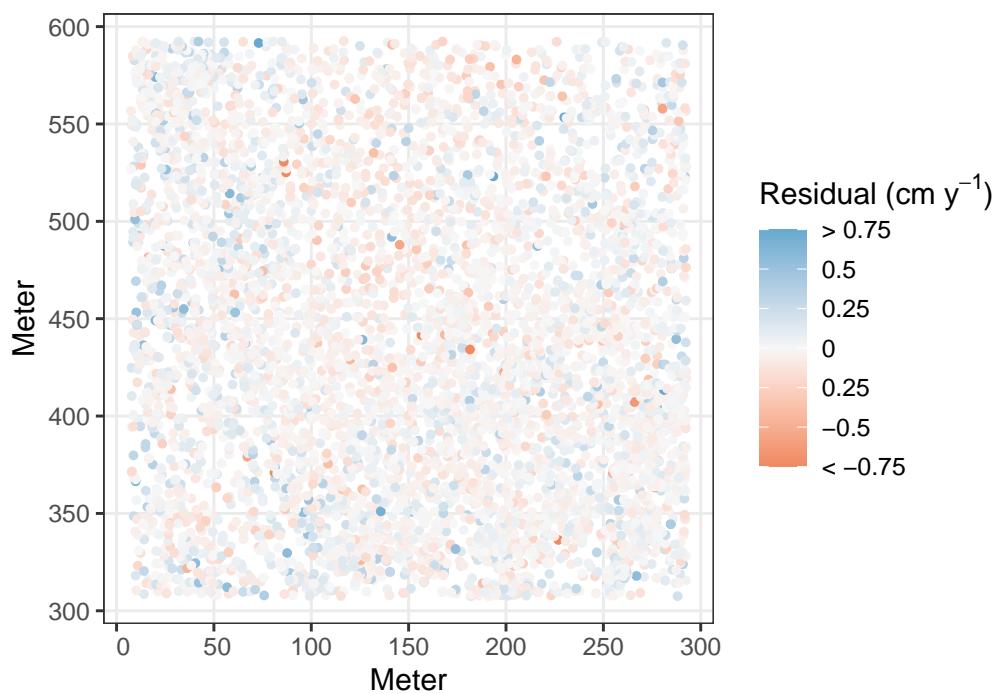


Figure 7: Spatial distribution of residuals for model applied to SCBI data.

239 **3.6 Run spatial cross-validation**

240 The model fits and predictions in Section 3.5 all suffer from a common failing: they use
241 the same data to both fit the model and to assess the model's performance using the
242 RMSE. As argued by Roberts et al. (2017), this can lead to overly optimistic assessments
243 of model quality as the models can be overfit, in particular in situations where spatial-
244 autocorrelation is present. To mitigate the effects of such overfitting, we use a spatially
245 block cross-validation algorithm implemented in the `run_cv()`. This function at its core
246 uses the same model fitting implemented in the `fit_bayesian_model()` function, however
247 trains the model on $k - 1$ spatial folds of the train and returns fitted values for the test
248 data. Recall that the spatial blocking scheme was encoded in Section 3.3.

249 **3.6.1 Big Woods**

250 Applying this spatially cross-validated model fit yields an RMSE is higher than that when
251 the model is fit without cross validation. In other words, our model fits in 3.5 were overly
252 optimistic in the model's fitting power, whereas a cross-validated results yield an estimate
253 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
```

254 **3.6.2 SCBI**

255 Observe once again that this RMSE is much higher than that for the above SCBI model
256 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

```

257 3.7 Visualize posterior distributions

258 Lastly, we return to the model fits from Section 3.5 and present tools to visually explore
 259 the posterior distributions of all parameters in our model. There are two main groups of
 260 parameters to consider. The β coefficients tell us about how fast each species grows and
 261 how this depends on DBH while the full matrix of λ values describe the competitive effects
 262 between pairs of species. There is a rich literature on this matrix (cite).

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

268 3.7.1 Big Woods

269 Here we re-run the model fit to the Big Woods data from Section 3.5, but this time use “fam-
 270 ily” as the group for comparison which has. This makes the posterior distributions easier to
 271 follow. Also, surprisingly, grouping by family performed just as well as grouping by species
 272 Allen & Kim (2020). First we re-run `create_focal_vs_comp()` and `fit_bayesian_model()`
 273 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)
```

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

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

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

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

281 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
```

282 Finally Figure 10 `lambda` gives the competition coefficients λ .

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

283 3.7.2 SCBI

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

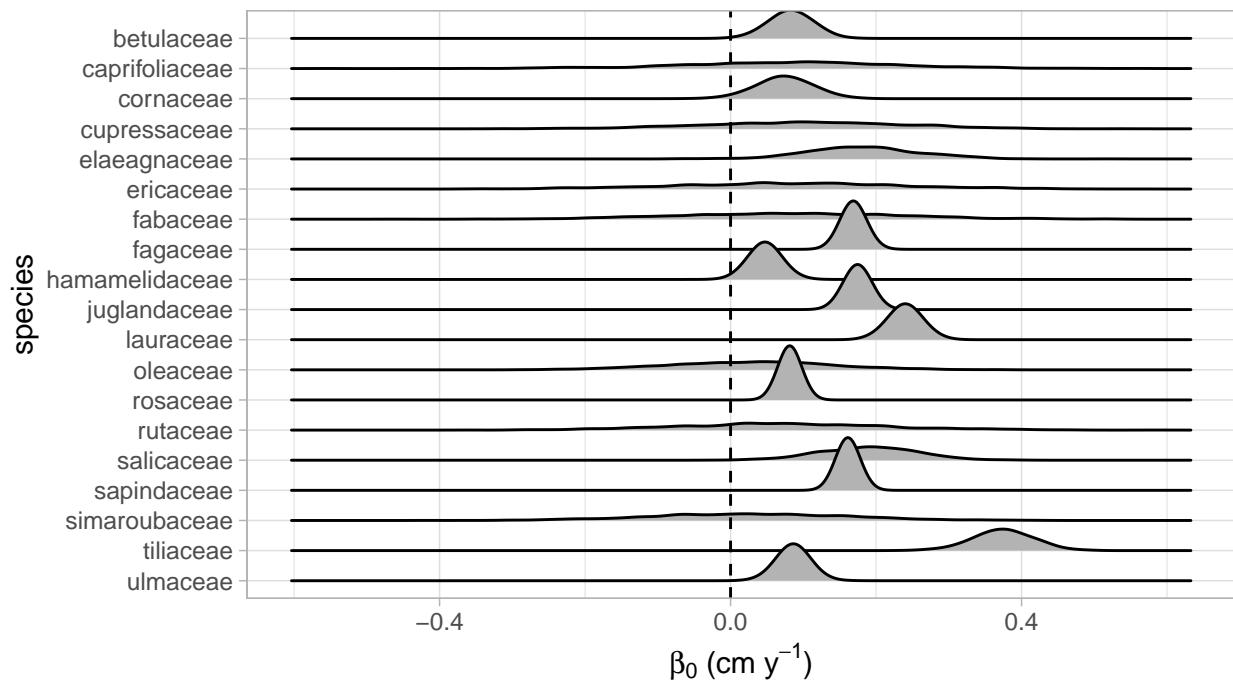


Figure 8: Posterior distribution of beta0.

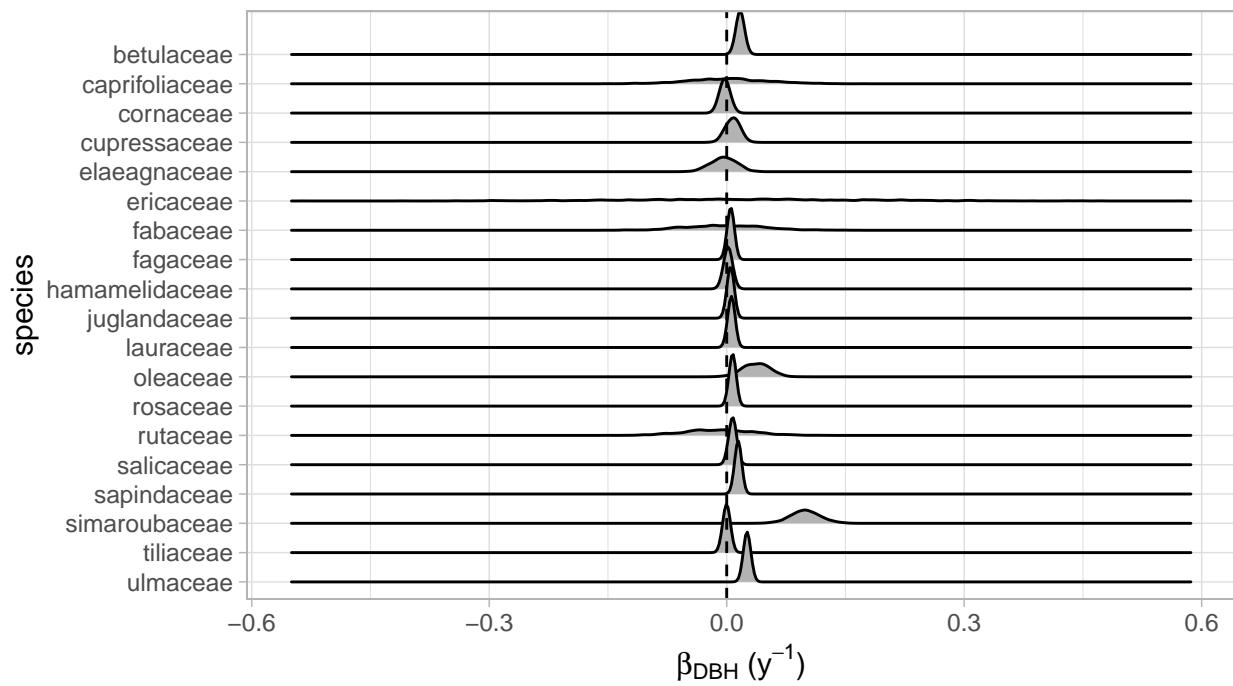


Figure 9: Posterior distribution of betadbh.

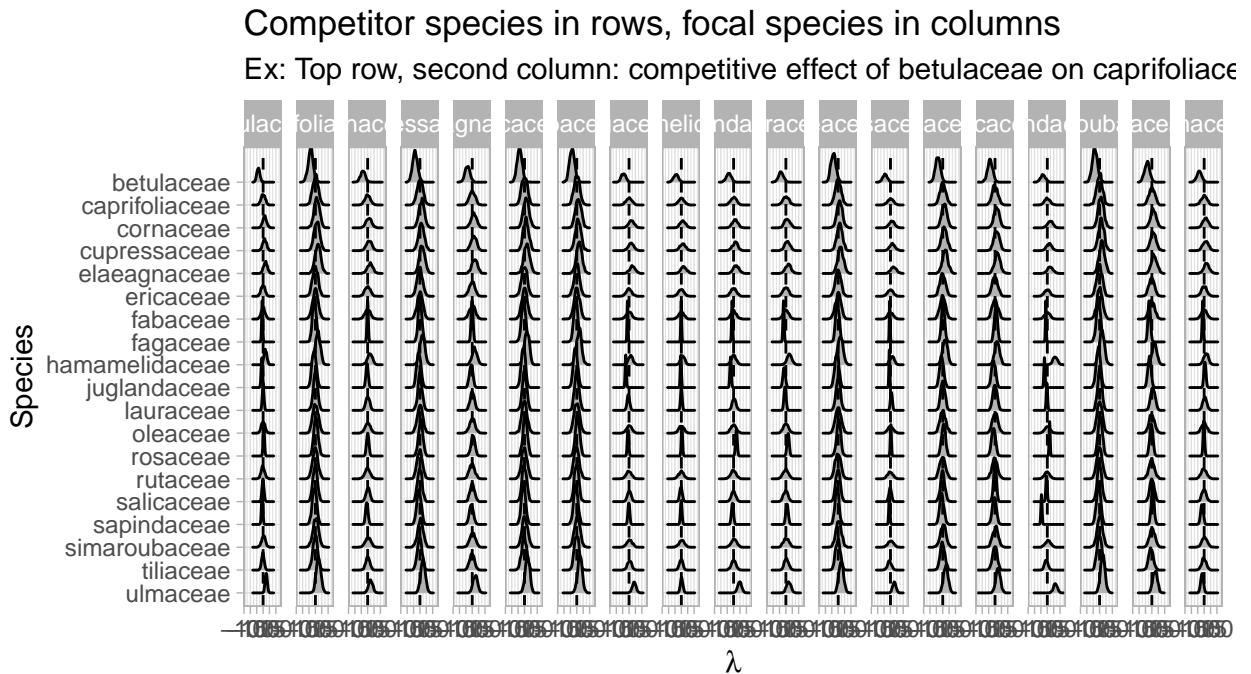


Figure 10: Posterior distribution of lambda's for Big Woods.

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

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

286 Add explanation here.

287 HEY BERT PICK IT UP HERE

288 4 Discussion

289 5 Acknowledgments

290 References

291 Allen, D., Dick, C., Burnham, R. J., Perfecto, I. & Vandermeer, J. (2020), 'The michigan big
292 woods research plot at the edwin s. george, pinckney, mi, usa', *Miscellaneous Publications*

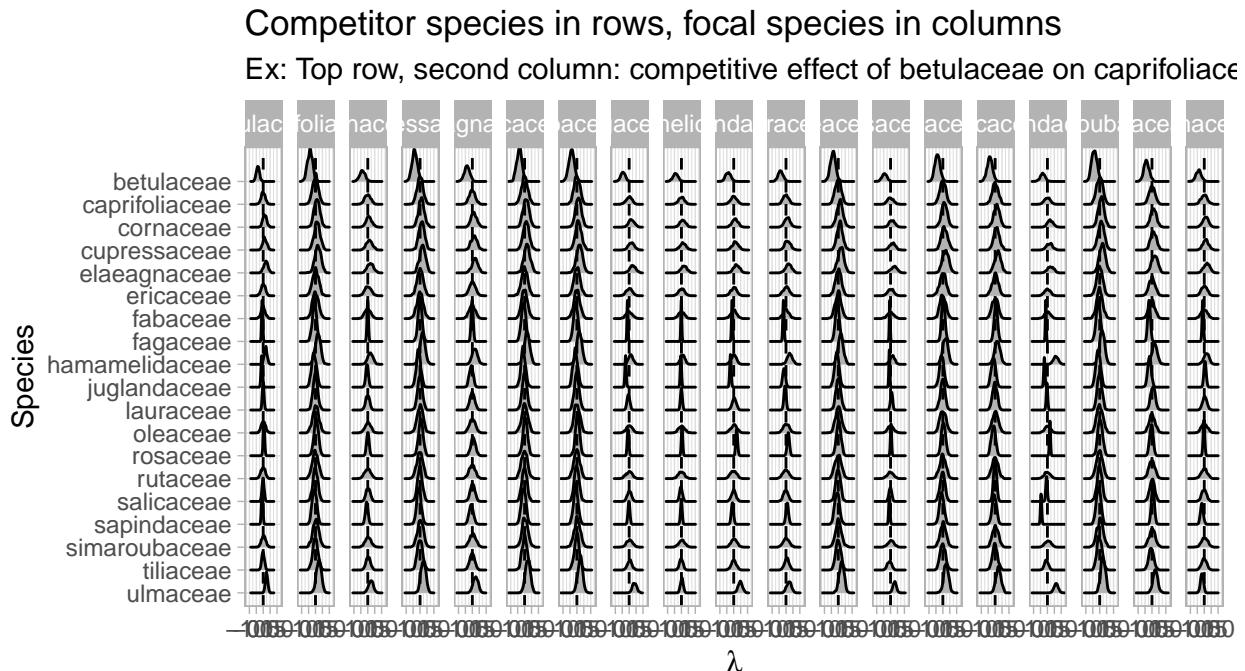


Figure 11: Posterior distribution of lambda's for SCBI.

293 *of the Museum of Zoology, University of Michigan* **207**.

294 **URL:** <http://hdl.handle.net/2027.42/156251>

295 Allen, D. & Kim, A. Y. (2020), ‘A permutation test and spatial cross-validation approach
296 to assess models of interspecific competition between trees’, *PLOS ONE* **15**(3), e0229930.

297 Publisher: Public Library of Science.

298 **URL:** <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0229930>

299 Anderson-Teixeira, K. J., Davies, S. J., Bennett, A. C., Gonzalez-Akre, E. B., Muller-
300 Landau, H. C., Wright, S. J., Salim, K. A., Zambrano, A. M. A., Alonso, A., Baltzer,
301 J. L., Basset, Y., Bourg, N. A., Broadbent, E. N., Brockelman, W. Y., Bunyavejchewin,
302 S., Burslem, D. F. R. P., Butt, N., Cao, M., Cardenas, D., Chuyong, G. B., Clay, K.,
303 Cordell, S., Dattaraja, H. S., Deng, X., Detto, M., Du, X., Duque, A., Erikson, D. L.,
304 Ewango, C. E. N., Fischer, G. A., Fletcher, C., Foster, R. B., Giardina, C. P., Gilbert,
305 G. S., Gunatilleke, N., Gunatilleke, S., Hao, Z., Hargrove, W. W., Hart, T. B., Hau, B.
306 C. H., He, F., Hoffman, F. M., Howe, R. W., Hubbell, S. P., Inman-Narahari, F. M.,
307 Jansen, P. A., Jiang, M., Johnson, D. J., Kanzaki, M., Kassim, A. R., Kenfack, D.,

308 Kibet, S., Kinnaird, M. F., Korte, L., Kral, K., Kumar, J., Larson, A. J., Li, Y., Li, X.,
309 Liu, S., Lum, S. K. Y., Lutz, J. A., Ma, K., Maddalena, D. M., Makana, J.-R., Malhi,
310 Y., Marthews, T., Serudin, R. M., McMahon, S. M., McShea, W. J., Memiaghe, H. R.,
311 Mi, X., Mizuno, T., Morecroft, M., Myers, J. A., Novotny, V., Oliveira, A. A. d., Ong,
312 P. S., Orwig, D. A., Ostertag, R., Ouden, J. d., Parker, G. G., Phillips, R. P., Sack, L.,
313 Sainge, M. N., Sang, W., Sri-ngernyuang, K., Sukumar, R., Sun, I.-F., Sungpalee, W.,
314 Suresh, H. S., Tan, S., Thomas, S. C., Thomas, D. W., Thompson, J., Turner, B. L.,
315 Uriarte, M., Valencia, R., Vallejo, M. I., Vicentini, A., Vrška, T., Wang, X., Wang, X.,
316 Weiblen, G., Wolf, A., Xu, H., Yap, S. & Zimmerman, J. (2015), 'CTFS-ForestGEO: a
317 worldwide network monitoring forests in an era of global change', *Global Change Biology*
318 **21**(2), 528–549.

319 **URL:** <http://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.12712>

320 Bourg, N. A., McShea, W. J., Thompson, J. R., McGarvey, J. C. & Shen, X. (2013), 'Initial
321 census, woody seedling, seed rain, and stand structure data for the SCBI SIGEO Large
322 Forest Dynamics Plot', *Ecology* **94**(9), 2111–2112.

323 **URL:** <http://esajournals.onlinelibrary.wiley.com/doi/abs/10.1890/13-0010.1>

324 Canham, C. D., LePage, P. T. & Coates, K. D. (2004), 'A neighborhood analysis of canopy
325 tree competition: effects of shading versus crowding', *Canadian Journal of Forest Re-*
326 *search* **34**(4). Publisher: NRC Research Press Ottawa, Canada.

327 **URL:** <https://cdnsciencepub.com/doi/abs/10.1139/x03-232>

328 Canham, C. D., Papaik, M. J., Uriarte, M., McWilliams, W. H., Jenkins, J. C.
329 & Twery, M. J. (2006), 'Neighborhood Analyses Of Canopy Tree Competi-
330 tion Along Environmental Gradients In New England Forests', *Ecological Applica-*
331 *tions* **16**(2), 540–554. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1890/1051-0761%282006%29016%5B0540%3ANAOC%5D2.0.CO%3B2>.

333 Gonzalez-Akre, E., McGregor, I., Anderson-Teixeira, K., Dow, C., Herrmann, V., Terrell,
334 A., Kim, A. Y., Vinod, N. & Helcoski, R. (2020), 'SCBI-ForestGEO/SCBI-ForestGEO-
335 Data: 2020 update'.

336 **URL:** <https://doi.org/10.5281/zenodo.4041595>

- 337 Pebesma, E. (2018), ‘Simple Features for R: Standardized Support for Spatial Vector Data’,
338 *The R Journal* **10**(1), 439–446.
339 **URL:** <https://journal.r-project.org/archive/2018/RJ-2018-009/index.html>
- 340 Pohjankukka, J., Pahikkala, T., Nevalainen, P. & Heikkonen, J. (2017), ‘Estimating the
341 prediction performance of spatial models via spatial k-fold cross validation’, *International
342 Journal of Geographical Information Science* **31**(10), 2001–2019.
- 343 Roberts, D. R., Bahn, V., Ciuti, S., Boyce, M. S., Elith, J., Guillera-Arroita, G., Hauen-
344 stein, S., Lahoz-Monfort, J. J., Schröder, B., Thuiller, W., Warton, D. I., Wintle, B. A.,
345 Hartig, F. & Dormann, C. F. (2017), ‘Cross-validation strategies for data with temporal,
346 spatial, hierarchical, or phylogenetic structure’, *Ecography* **40**(8), 913–929.
347 **URL:** <http://onlinelibrary.wiley.com/doi/abs/10.1111/ecog.02881>
- 348 Smith, W. B. (2002), ‘Forest inventory and analysis: a national inventory and monitoring
349 program’, *Environmental pollution* **116**, S233–S242.
- 350 Tatsumi, S., Owari, T., Ohkawa, A. & Nakagawa, Y. (2013), ‘Bayesian modeling of neigh-
351 borhood competition in uneven-aged mixed-species stands’, *Formath* **12**, 191–209.
- 352 Uriarte, M., Condit, R., Canham, C. D. & Hubbell, S. P. (2004), ‘A spa-
353 tially explicit model of sapling growth in a tropical forest: does the iden-
354 tity of neighbours matter?’, *Journal of Ecology* **92**(2), 348–360. _eprint:
355 <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.0022-0477.2004.00867.x>.
356 **URL:** <http://besjournals.onlinelibrary.wiley.com/doi/abs/10.1111/j.0022-0477.2004.00867.x>
- 358 Valavi, R., Elith, J., Lahoz-Monfort, J. J. & Guillera-Arroita, G. (2019), ‘blockCV: An
359 r package for generating spatially or environmentally separated folds for k-fold cross-
360 validation of species distribution models’, *Methods in Ecology and Evolution* **10**(2), 225–
361 232. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/2041-210X.13107>.
362 **URL:** <http://besjournals.onlinelibrary.wiley.com/doi/abs/10.1111/2041-210X.13107>
- 363 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., Grole-
364 mund, G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T. L., Miller, E., Bache,

³⁶⁵ S. M., Müller, K., Ooms, J., Robinson, D., Seidel, D. P., Spinu, V., Takahashi, K.,
³⁶⁶ Vaughan, D., Wilke, C., Woo, K. & Yutani, H. (2019), ‘Welcome to the Tidyverse’,
³⁶⁷ *Journal of Open Source Software* 4(43), 1686.
³⁶⁸ **URL:** <https://joss.theoj.org/papers/10.21105/joss.01686>