

1 Predicting biodiversity dynamics in response to environmental
2 change

3 Can we do it? A report from assess.sim.basic.R

4 Ryan Batt

5 2015-08-23

6 **Abstract**

7 “Lorem ipsum dolor sit amet, consectetur adipisicing elit, sed do eiusmod tempor incididunt ut labore
8 et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut
9 aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum
10 dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia
11 deserunt mollit anim id est laborum.”

¹² Contents

¹³	Introduction	4
¹⁴	Overview	4
¹⁵	The Simulation	4
¹⁶	Multispecies Occupancy Models (MSOMs)	5
¹⁷	Conventions and Settings	7
¹⁸	Dimension Conventions	7
¹⁹	Settings	8
²⁰	Species Richness	11
²¹	Definition of species richness	11
²²	Regional Richness	11
²³	Site Specific Richness (N_{site})	13
²⁴	Occupancy Probability, ψ	16
²⁵	Definition of ψ	16
²⁶	Scatter Plot of Aggregated ψ	16
²⁷	Scatter Plot of $\hat{\psi}$ vs ψ_{true} , split by year and replicate	17
²⁸	Occupancy Response Curves	17
²⁹	Probability of Detection, p	21
³⁰	Definition of p	21
³¹	Demo: Effect of MSOM Hierarchy on p	21
³²	Scatter Plot of \hat{p} vs p_{true}	23
³³	Scatter Plot of \hat{p} vs p_{true} , split by year and replicate	23
³⁴	Assessment with Mixed Effects Models	25
³⁵	Describe Motivation for Mixed Effects Models	25
³⁶	Example LMER Analysis for ψ	25

37	Report Generation Notes	27
38	R Session Information	27
39	Date Document Last Compiled	27

40

41 Introduction

42 Overview

43 As water temperatures change, species may shift the size and location of their geographical ranges, bearing
44 consequences for the food webs and economies linked to those species. However, species don't always respond
45 similarly to shifting temperatures (different thermal tolerances, e.g.), which means that changing temperature
46 may remix the composition and diversity of ecological communities.

47 The biological, spatial, and temporal scale of community diversity shifting in response to climate is massive.
48 A functional definition of a community may consist of 100's or 1000's of species, each of which may be
49 shifting its range at a scale of decades and 100's kilometers. As a result, we need statistical methods for
50 estimating biodiversity that don't rely on heavy replication and that make efficient use of available data.
51 Enter the superstars: on the data side the trawl data set has amazing spatiotemporal and taxonomic extent
52 and resolution; on the statistical side multispecies occupancy models (MSOM) are hierarchical state space
53 models that are designed to estimate species richness and don't require consistent or extensive "replication".
54 Although they're superstars, even these data and models have their limitations and pitfalls.

55 Can we estimate the dynamics of species richness from trawl data using an MSOM? It's a hard question to
56 answer because we can never know the "truth" for sure, but we can get an idea of how reliable our analysis
57 is by simulating fake data, for which we know true values because we created them. The trawl data set is
58 generated by two distinct processes: Nature's data generating process (NDGP), and the process by which
59 humans observe the result of NDGP. So we ask: to what extent is the accuracy of estimates from an MSOM
60 dependent on characteristics of NDGP, and in particular, the way in which we observe the result of NDGP?
61 The strategy for answering this question is to simulate fake data where we approximate Nature but gain
62 knowledge of "truth", "observe" the results of the true process, then try to recover the true species richness
63 from these simulated data.

64 The Simulation

65 The goal of this simulation was to use a very basic process to generate presences and absences of species in
66 space and time. In this version of the simulation, there is no explicit connection between years (they are
67 independent). There is a modest spatial connection, because in the simulation an environmental variable
68 determines habitat suitability. I think of this environmental variable as temperature, and I filled a grid with
69 temperatures that ranged from the coldest at the top of the grid (north) and the warmest at the bottom
70 (south) and added random variation among columns in the same row (among longitudes at the same latitude).

71 One level of the simulation mimics NDGP. In this level, NDGP is best characterized by ψ , which is the product
72 of a temperature and species' response curves. I.e., temperatures were used to determine the suitability of
73 each grid cell to each simulated species. This suitability is known as ψ throughout this document.

74 A second level of the simulation mimics human observation of NDGP — what we do when we collect data.
75 This process was simulated by assigning each species has a unique probability of being observed or "detected"
76 (this variable is p). The observation process gets several attempts at observing a given species in a given grid
77 cell; think of this as subdividing each site into subsites, and when you visit each subsite you have probability

78 p of observing a particular species (each species has its own p). Depending on the settings used in the analysis
79 that this document summarizes, the maximum number of subsites can vary, as can the number of subsites per
80 site (OK, fine; the maximum number of subsites in this version is 4, the number of subsites per site varied
81 between 1 and 4, and overall 50% of total possible subsites were sampled).

82 As previously mentioned, the simulation included “time”. In this basic version, not much changes between
83 the “years” for the true process (temperature doesn’t change, nor do the response curves), but the mean of p
84 does change. In a given year, the entire community has an overall mean probability of being detected, and
85 each species randomly deviates from that mean.

86 The simulation also has replicates. To understand the replicates, it needs to be clear that even when a
87 parameter in the simulation does not change, the outcome can change. The replicates hold the realization
88 of the simulated NDGP constant, and draw new realizations of the observation process. I.e., both ψ and
89 p are constant among replicates, and the binary *outcome* of ψ is also held constant, but the outcome for p
90 can change. Furthermore, although each replicate has same values of p (both the mean p and each species’
91 individualized random draw from that distribution), each replicate switches which year is associated with
92 which p ’s. In this way we can observe each outcome of Nature’s data generating process under a series of
93 settings for the human observation process.

94 Multispecies Occupancy Models (MSOMs)

95 Multispecies occupancy models are Bayesian statespace hierarchical models. They distinguish between truth
96 and observation of the truth, and many parameters share a common “parent” distribution. They are very
97 flexible models, and can be adapted to include new types of processes. The MSOM being used here is a
98 relatively simple version of these models. It predicts the probability of each species existing in a grid cell from
99 a logistic regression equation that uses a second-order polynomial of the environmental variable as a covariate.
100 The parameters in this level of the model are hierarchical, with species having their own paramter values,
101 but these individual parameters are not wholly independent in the sense that they share a common parent
102 distribution, which sort of acts to both limit how different they can be and to inform one another. The model
103 also has an observation level, which only has a hierarchical intercept (just a mean) as a predictor variable.
104 The MSOM makes guesses of the true state of the system (whether a species is actually present or not). It
105 then makes guesses at how the observation of that true state might turn out, which is effectively a prediction
106 of what our data will be. The Bayesian model fitting process then uses this comparison of the observed data
107 to the estimate of the observation to tweak the parameters in the MSOM. This process is repeated until the
108 choice of paramters boils down to what is essentially the posterior distribution of the estimated parameters.
109 Right now the MSOM model is fit separately to each year and each replicate. So the model never gets to see
110 multiple years or multiple replicates at the same time. Furthermore, when referring to a parameter value
111 fitted in the MSOM, it is implied that it can be subscripted with time or replicate (because all years and
112 replicates are fit independently).
113 The parameters in the logistic regression that predicts the value of ψ vary among species, although ψ itself
114 varies among species and space, because the regression parameters (subscripted by species) are multiplied by
115 the environmental variable (subscripted by space). More or less, it can be said that, for a given species, ψ

₁₁₆ varies among space because of the environmental variable, and in a given location it varies among species
₁₁₇ because of the regression parameters.

₁₁₈

119 Conventions and Settings

120 In this section I outline the subscripting and notation used in the MSOM analysis and for the simulation. I
121 also outline various settings (number of species simulated, replicates, etc.). Most of the numbers you see
122 (and some of the text) is dynamically generated based on the code that produced the statistics and figures.
123 Therefore, you can refer back to these sections to see what settings may have changed since the last version
124 of this document.

125 **Note:** *I've often found myself having to get creative with subscripts and superscripts. I've tried to be clear an*
126 *consistent, but small inconsistencies likely exist, so don't be confused by them. For example, if you see $\max(Z)$*
127 *and Z_{\max} in two different sections, they are probably referring to the same thing. If you see something*
128 *confusing, let me know (preferably by (creating an issue on GitHub)[<https://github.com/rBatt/trawl/issues>]),*
129 *and I'll fix it.*

130 Dimension Conventions

131 Summary

132 1. Site ($j = 1, 2, \dots, j_{\max} = 9 \times 9 = 81$)

- 133 • Sites are unique combinations of latitude and longitude
- 135 • The spatial arrangement of sites is *not* arbitrary, although importance depends on settings (see
136 `dynamic` below)
- 138 • The environmental variable X varies among sites (and years, below)

139 2. Sub-sites ($k = 1, 2, \dots$)

- 140 • Sub-sites are only relevant to the “observation” process
- 141 • Each site has the same number of possible sub-sites, but the number of sub-sites observed can vary
- 142 • In this simulation, $k_{\max} = 4$, $k_{\min}^{\text{observed}} = 1$, and $k_{\max}^{\text{observed}} = 4$
- 144 • Substrata are primarily useful for determining p , the **detection probability**

146 3. Species ($i = 1, 2, \dots i_{\max} = R = 30$)

- 147 • Does not include “augmented” species
- 148 • For this MSOM analysis, the species array was padded with 10 0’s

149 4. Time ($t = 1, 2, \dots 4$)

- 150 • Time is primarily used to vary the parameters controlling the “true” process
- 151 • When those parameters don’t change, time provides independent*realizations of the same “true”
- 152 process

153 — *Note: only when `dynamic=FALSE` in `sim.spp.proc`

154 5. Replicates ($r = 4$)

- 155 • Replicates are *simulated* repeated human observations of the same *realization* of the “true” process
156 at Time t
157 • Replicates are used to vary the parameters that control the “observation” process
158 • When those parameters don’t change, each replicate provides an independent* realization of the
159 same “observation” process

160 **In Code**

161 The MSOM analyzes each year t -replicate r combination independently. Parameters subscripted by these
162 dimensions are derived from separate analyses.

163 In my code, I’ve tried to be consistent in my use of these indices to describe arrays, matrices, and rasters.
164 Rows are dimension 1, columns dimension 2, etc. The precedent of subscripted dimensions follows the numeric
165 ordering of the list above. E.g., in the matrix $X_{i,t}$ each row will refer to a different species, and each column
166 a different year (note that site j is skipped, so species i is “promoted” to dimension 1, the row.). By default, R
167 fills matrices and arrays by column, whereas the **raster** package fills them by row. In most cases where an R
168 object needs to split sites into the lat/ lot components, I make use of the **raster** package. Therefore, the
169 numbering of the sites proceeds row-wise, where each site is numbered according to the order in which it is
170 filled, as in this 2×3 matrix: $J = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix}$

171 Note that even though this matrix is numbered row-wise, it is still indexed as $J_{row,column}$, such that $J_{1,2} = 2$.
172 As mentioned previously, this information is primarily important for understanding the code involved with
173 this project, and in most cases it is not crucial to be explicitly aware of the spatial arrangement of sites.

174

175 **Settings**

176 **Simulation Settings**

177 I created a class called "spp", which has methods for **print()**. The **Dimensions** are the number of sites, the
178 number of species, then the number of years.

179 Also **printed** are some richness summary statistics. **All cells** refers to the collective richness over all j
180 taken together. The meaning of **One cell** differs slightly between the true and observed printouts: in the
181 true printout the richness is of a particular site (j), and in the observed printout it is of a particular sub-site
182 (k).

183 ## Dimensions: 81, 30, 4
184 ## grid.h = 9
185 ## grid.w = 9
186 ## grid.t = 4
187 ##
188 ## Number Species Possible (ns):

```

189 ## 30
190 ## Total Species Richness:
191 ## 30
192 ## Total Observed Species Richness:
193 ## 30
194 ##
195 ## Annual Species Richness:
196 ##      Min. 1st Qu. Median Mean 3rd Qu. Max.
197 ## All cells 29      29     29 29.25 29.25   30
198 ## One cell  5       10    14 13.63 17.00   24
199 ##
200 ##
201 ## Observed Annual Species Richness:
202 ##      Min. 1st Qu. Median Mean 3rd Qu. Max.
203 ## All cells 27      27.75 28.5 28.250 29     29
204 ## One cell  0       0.00   0.0  4.008   7     23

```

205 In the MSOM, detectability (p_i) is determined in the form of a logistic regression, which currently only
206 has an intercept (v_0) as predictor (so just a mean). That intercept varies among species (i.e., $v_{0,i}$), and
207 that variation is generated by drawing each individual species's intercept ($v_{0,i}$) from a parent distribution:
208 $v_{0,i} \sim \mathcal{N}(\mu_{v_0}, \sigma_{v_0}^2)$. See [section about \$p\$](#) for more info.

year	mu.v0	sigma.v0
1	-2	2
2	0	2
3	2	2
4	4	2

209

210 Settings for JAGS & MSOM

nChains	nIter	n0s	nSamples
3	50000	10	500

211 In the table above, `nChains`, `nIter`, and `nSamples` are all variables that are strictly pertinent to the Bayesian
212 analysis carried out in JAGS. The `n0s` value refers to the the degree of “data augmentation”. In this process,
213 you add extra species to the data set, and say that they were never observed. For our purposes, this is
214 employed for purely technical reasons, although it can be used to extra further inferences about species
215 richness.
216 The posterior samples from JAGS consist of 500 samples (above). However, to save memory and hard drive

217 space, I have often only saved measures of central tendency for each of these. In this assessment, I have
218 performed all calculations on the **centralT=median** of the posterior samples.

219

220 **Species Richness**

221 **Definition of species richness**

222 Species richness is the number of different species, or more generically, unique taxa. The point is moot in the
223 simulation study, and in the empirical trawl data it refers to species.

224 Estimates of richness (R) can be made spatially or temporally explicit (or neither, or both). In the following
225 figures, different levels of aggregation are performed – for most figures R is split by year (this is true for
226 all figures but the Boxplot Figure). The Time Series of Richness Figure emphasizes temporal dynamics
227 and keeps replicates separated, but aggregates over space (the j sites). The Nsite Scatter Figure doesn't
228 aggregate over space or time, but it does aggregate over “replicate” observations; importantly, while the
229 figure does present any spatial aggregation, it does not retain the spatial relationship (you can't tell which
230 sites are next to others). The final two figures of the section (Heatmap of of Richness Figure) are similar to
231 the previous figure, except that spatial relationship among points is retained via a heatmap representation.

232 None of these estimates of richness include the 10 species that were part of the “data augmented”/ “adding
233 0's” process. Richness values can either be true (true simulated NDGP; R^{true}), observed (true simulated
234 human observation of NDGP; R^{obs}), or MSOM estimates of one of those two (\hat{R}^{true} or \hat{R}^{obs}).

235

236 **Regional Richness**

237 These estimates of species richness only distinguish between replicates and years. They do not contain any
238 site-specific information.

239 **Richness Boxplots**

240 With the boxplots we're mostly looking to see if the estimates of richness vary with the mean probability of
241 detection, p . In the empirical data, we know that taxonomic identification changed over time (it improved;
242 generally, more species were ID'd in later years). We also suspect that gear might change, which affects
243 the probability of observing a species. The “Average Detection Probability” category in the boxplots is the
244 cross-species average of p (which with large sample size approach the hyperparameter p_μ).

245

246 **Richness Time Series**

247 Text explanation goes here

248 Need explanations for how each panel was calculated.

249 1. R^{true} is straightforward

250

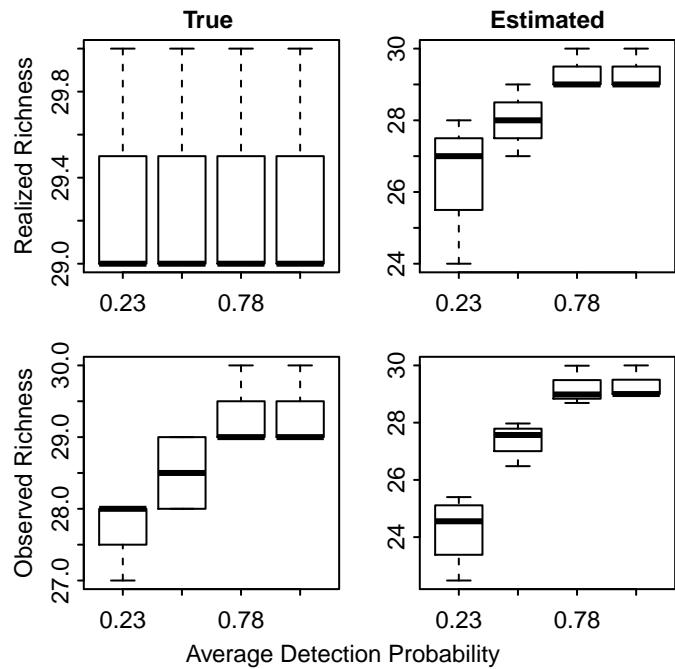


Figure 1: Boxplots of species richness. Numeric groupings indicate the average value of p across species during a given year–replicate combination. The panels in the left column are the true simulated values, and the panels on the right are the corresponding MSOM estimates. The top row indicates the latent realized species richness or MSOM estimates of the richness. The bottom row’s panels are the simulated observed values of richness (the response variable in the MSOM) and the MSOM estimates of the observed values.

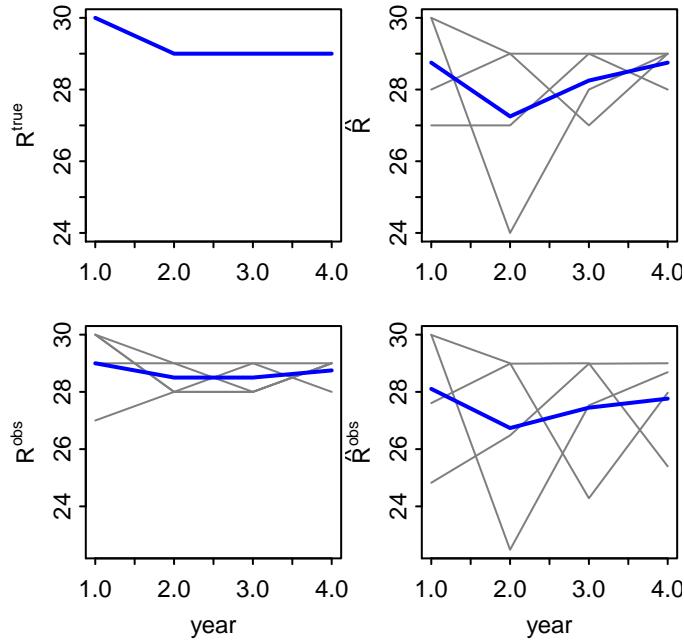


Figure 2: Time series of richness. Gray lines are individual replicates. Blue lines are averages. Note that detection probabilities ($p_{i,t,r}$, see [simulation settings above](#), as well as [definition of \$p\$ below](#)) change over time, and their temporal ordering differs among replicates.

251 2. \hat{R} is from $\sum_{i=1}^{R^{true}} \max(\hat{Z}_{i,\nabla j \in J})$; and to be clear, \hat{R} does not include the “unobserved” species introduced
252 to the MSOM occurrence matrix (Y)

253

254 **Site Specific Richness (Nsite)**

255 **Scatter Plots of Nsite Split by Year**

256

257 **Maps of Richness (space and time)**

258 Text explanation goes here

259

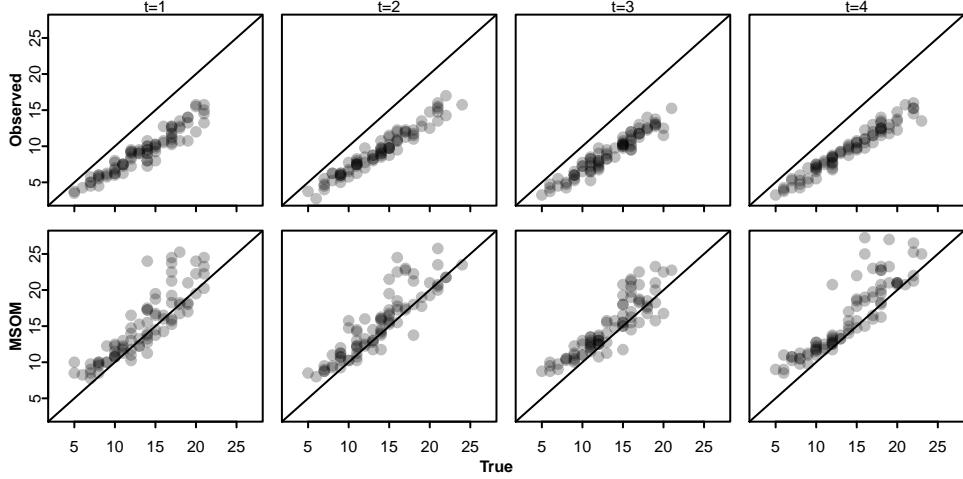


Figure 3: Site-specific richness (`Nsite`, N_j) from simulated observations (vertical axis, top row; N_j^{obs}) and from MSOM estimates (vertical axis, bottom row, \hat{N}_j) vs true site-specific richness (horizontal axis; N_j^*). The panel columns delineate the years of the simulation. Each point is site-specific species richness that has been averaged over the simulated replicate observations.

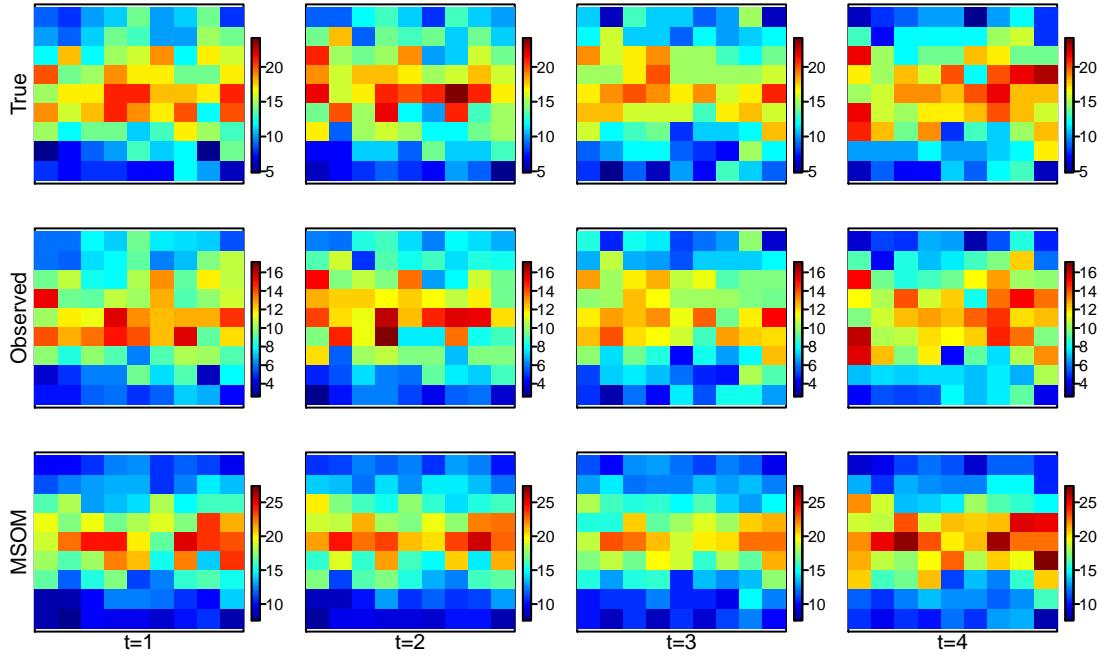


Figure 4: Maps of site- and year-specific species richness (`Nsite`) from the simulation of the True process (top row), simulation of the Observed process (middle row), and the MSOM estimates (bottom row). X-axis and Y-axis indicate position in 2 dimensional space; it is important to note that the environmental variable changes linearly across the y-axis, and randomly (and much less) across the x-axis. The different columns represent separate years. The environmental variable changes linearly among years (the rate of change is the same for all x-y locations). Colors indicate species richness (warm colors are higher richness than cool colors), averaged over the simulated replicate observations. Horizontal and vertical axes Each row of panels is scaled independently, columns within a row are scaled equally.

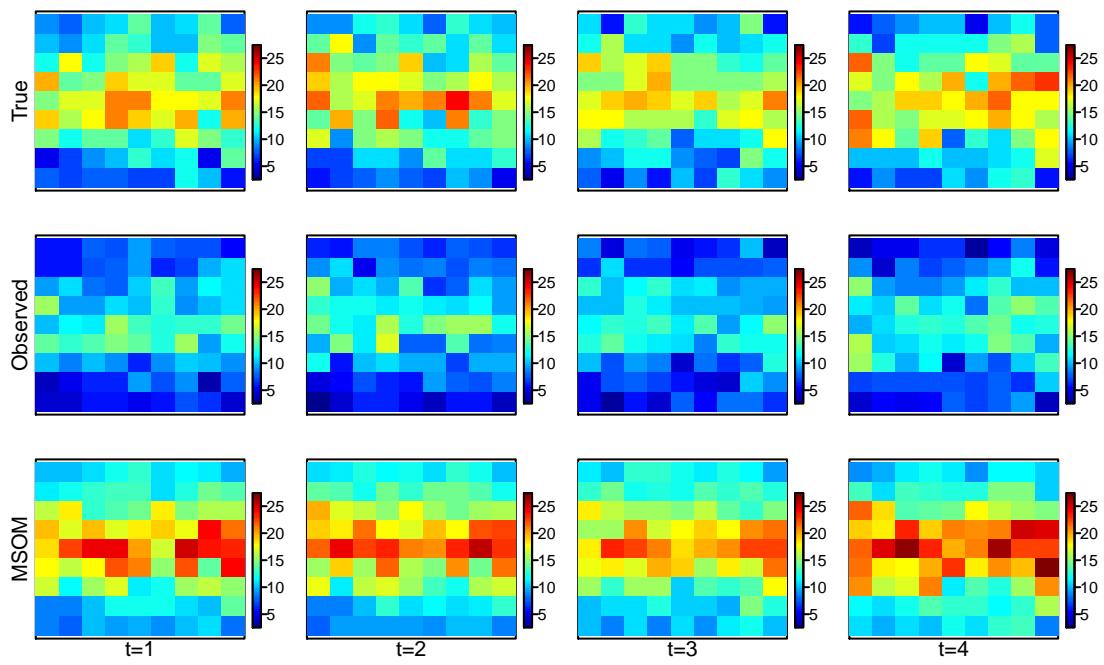


Figure 5: Same as previous figure, but all panels are on the same scale.

260 **Occupancy Probability, ψ**

261 **Definition of ψ**

- 262 Definition description goes here
- 263 Probably need to describe how it's generated in the simulation
- 264 As well as how it's estimated in the MSOM
- 265 In particular, important to point out that they may or may not match

266

267 **Scatter Plot of Aggregated ψ**

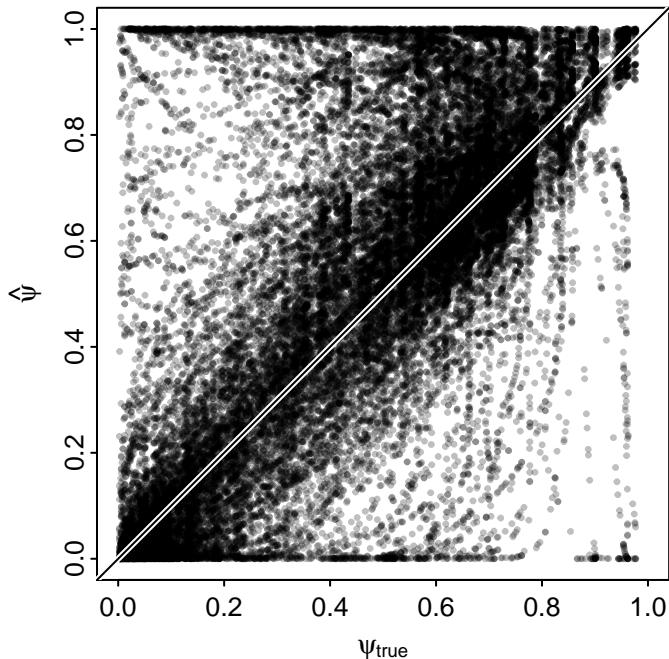


Figure 6: MSOM estimates of ψ ($\hat{\psi}$) vs. true values of ψ (ψ_{true}). Each point is a ψ value for a particular site-species-year-replicate. The white and black line is the 1:1 line.

- 268 In a general sense, the MSOM can distinguish between instances (sites/ years) when a species is likely to be present, and when it's not. However, in every simulation I've done (varying many parameters that aren't compared in this document), the scatter plot of ψ always makes it apparent that

271 1. There is a lot of variability around the 1:1 line

272

273 2. The residuals are not normal, and they are not independent

274 i. In general, I've found that $\hat{\psi}$ exhibits an upward bias, overestimating ψ^{true}

275

276 ii. Smoothly-curving excursions from the 1:1 line often prominent

277 These patterns are somewhat concerning. The curve-like sequence of residuals is probably a byproduct of
278 slightly incorrect estimates of the parameters in the logistic regression ($[a_0, a_1, a_2]$), resulting in estimated
279 **response curves** that deviate non-randomly from the true response curve. For a heuristic of how these
280 smooth excursions can occur, in R try something as simply as `d <- rnorm(100); plot(dnorm(d), dt(d,`
281 `1))` to see the relationship between the density estimate from the correct distribution and that from
282 the wrong distribution (the density is analogous to ψ); or for really crazy patterns, try `d <- rnorm(100);`
283 `plot(dnorm(d), do.call(approxfun, density(d)[c("x", "y")])(d))`. So the curves are explainable, but
284 I cannot explain the consistent overestimation; I could understand how underestimating detectability (p)
285 would result in overestimating ψ , but the MSOM appears to recover true p values rather well (e.g., see **P**
286 **Scatter Figure**), so that's not a satisfying explanation.

287 In the next section I drill into ψ a bit more to try and understand what causes the largest deviations from
288 true values.

289

290 Scatter Plot of $\hat{\psi}$ vs ψ_{true} , split by year and replicate

291 The estimates and true values of ψ are best correlated when p is high. When the average species has a low
292 chance of being detected (when p_μ is, say, 20%), the estimates of ψ are a mess.

293 *Note: what I refer to as p here is really just the probability that a species will be detected if an occupied site is
294 sampled, so the number of substrata sampled per site isn't reflected in p . In this simulation, 50% of substrata
295 were sampled, and while this doesn't influence p , it could add noise to its estimates.*

296

297 Occupancy Response Curves

Occupancy response curves are calculated as $\text{logit}(\psi_i) = \mathbf{X} \times \mathbf{a}_i$, where

$$\mathbf{X} = \begin{pmatrix} 1 & X_{min} & X_{min}^2 \\ \vdots & \vdots & \vdots \\ 1 & X_{max} & X_{min}^2 \end{pmatrix}; \mathbf{a}_i = \begin{pmatrix} a_{0,i} \\ a_{3,i} \\ a_{4,i} \end{pmatrix}$$

298 Therefore, these curves are tantamount to values of ψ , except that ψ generally pertains to a simulated,
299 observed, or true occupancy probability, whereas the occupancy probability in the response curves is calculated
300 over hypothetical conditions (i.e., over hypothetical values of the environmental gradient X).

301 True Occupancy Response Curves

302 In the response curve, the values of the environmental variable are an arbitrary gradient, and do not necessarily
303 correspond to what was observed in the simulated environment (although they are intended to cover the
304 same range).

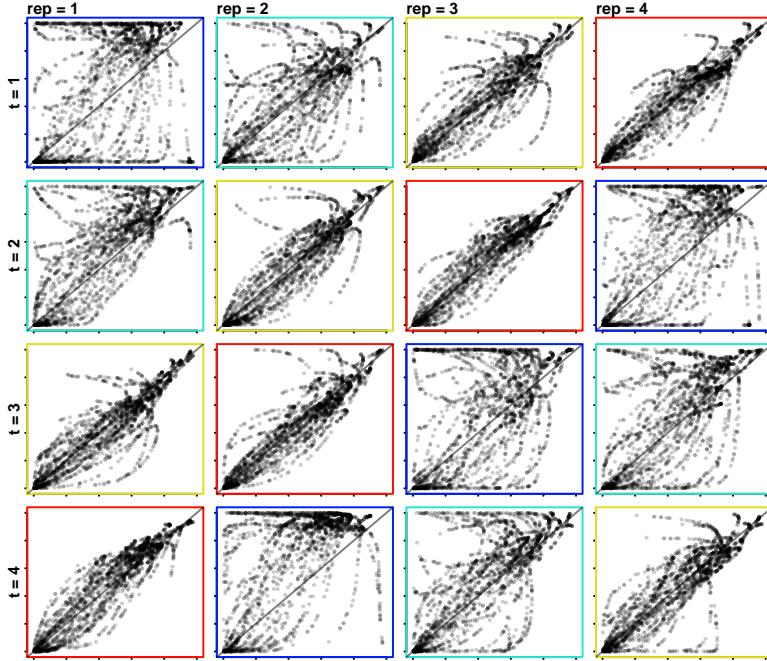


Figure 7: True (horizontal axes) and MSOM estimates (vertical axes) of occupancy probabilities ($\psi_{j,i,t,r}$) of species i occupying a location j . Years (t) are separated by rows, replicates (r) are separated by columns. The border color of each panel indicates the community-level mean probability of detection (p_μ ; where $p_i \sim \mathcal{N}(p_\mu, \sigma^2)$), with warm colors indicating high detectability, and cool colors low. The species-specific detectabilities are **not** re-randomized among replicates, but even when the probabilities associated with the observation process do not change, the outcome of the process can change. The year t of the simulated true process changes across the rows of panels, and the simulated replicate observation r changes across columns.

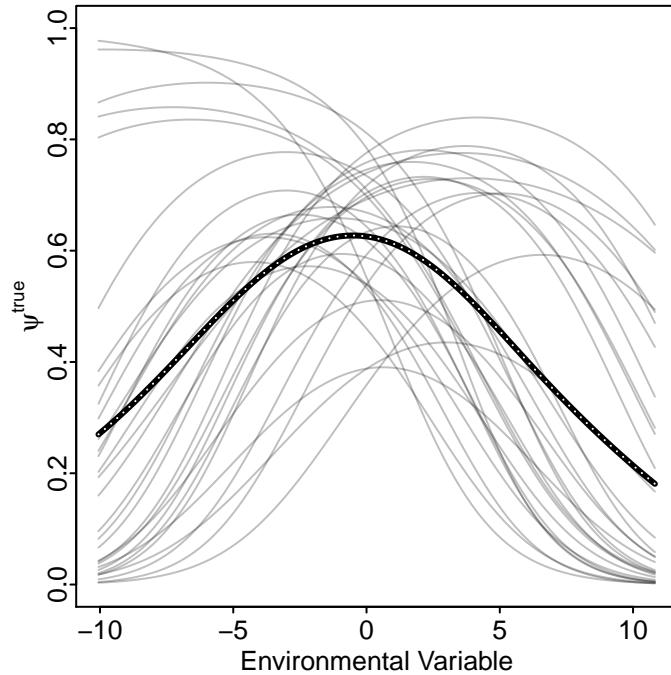


Figure 8: True simulated response curves. Vertical axis is the value of ψ^{true} , horizontal axis is the value of the environmental variable that, along with species-specific regression parameters, determines ψ^{true} . The thick line is the among-species mean value of ψ^{true} at a given value of the environmental variable.

305 **Estimated Occupancy Response Curves**

- 306 (min(X) = -10.1, and max(X) = 10.8)
 307 (red; $p_{max} = 0.93$)
 308 (blue; $p_{min} = 0.23$)

309

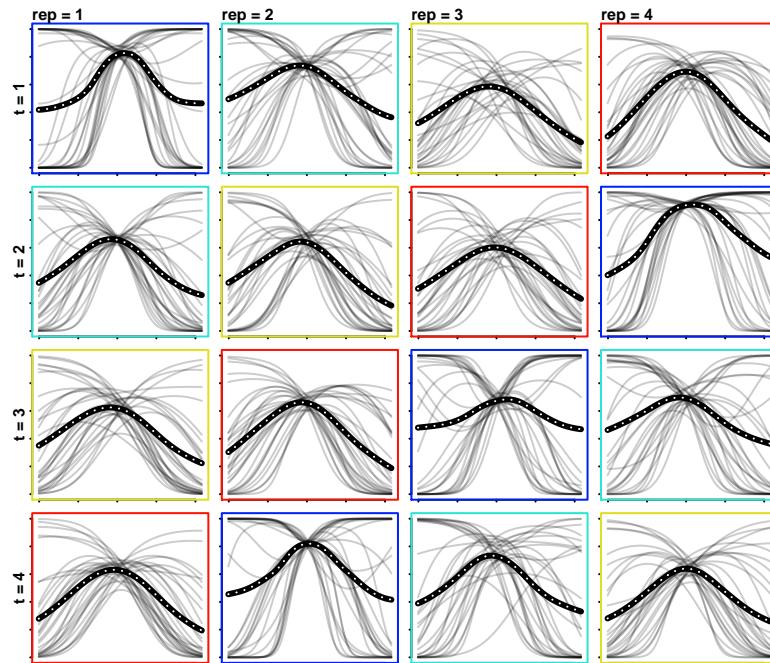


Figure 9: Response curves of species' probability of occupancy (ψ_i , vertical axis) across the full range of temperatures in the simulation. The color of the boxes around each panel refer to the among-species average of the probability of detection; warm colors indicate that the mean detection probability is high, whereas cool colors indicate that p was low. The year t of the simulated true process changes across the rows of panels, and the simulated replicate observation r changes across columns.

310 **Probability of Detection, p**

311 **Definition of p**

312 The probability of detection (p), is a species specific parameter in the MSOM model. The MSOM analyzes
313 all years (t) and replicates (r) separately, so I am going to leave those subscripts out of this description. In
314 the simulation, the probability of observing a species is a function of two independent factors:

- 315 1. The probability that site j is occupied by species i ; this is $\psi_{j,i}$
 - 316 • $\psi_{j,i}$ is a function of species-specific niche and an environmental variable that changes over space
317 and time
 - 319 • $Z_{j,i}$ is the species- and site-specific richness, which is a function of ψ (given that we're only talking
320 about species that are in the pool of possible species, determined by w_i)
- 321 2. A species-specific (i) chance of being identified (`taxChance`), given that it is present in a location that
322 was sampled (i.e., detectability does not reflect the probability of sampling a place); this detectability
323 parameter is p_i
 - 324 • Detectability changed between years.
 - 325 • In a given year, $\text{logit}(p_i) \sim \mathcal{N}(\sqrt{\mu}, \sigma^\epsilon)$. μ changed between years (taking on values of -2, 0, 2,
326 and 4), $\sigma^2 = 2$ in all years.
 - 327 • The value of p only changes between species (and years), but the observation process occurs at the
328 substratum (k) level. Thus, the parameter is really $p_{j,k,i}$, but for a given i , all $p_{j,k}$ are constant. I
329 represent this probability as p_i with the understanding that this value is repeated over space.
 - 330 • $Y_{j,i}$ is the observed version of $Z_{j,i}$.
 - 331 • $Y_{j,i} \sim \text{Bern}(p_i \times Z_{j,i})$.
 - 332 – Note: Because p is actually subscripted to k , the Y are also actually subscripted to k . Maybe
333 leaving these subscripts out is making things more confusing. I've only excluded them to
334 emphasize how parameters are estimated.
 - 335 • Our data about species presence/ absence correspond to $Y_{j,i}$. So it might be useful to think of the
336 MSOM as estimating $\hat{Y}_{j,i}$, which is compared to the observed data $Y_{j,i}^{obs}$.

337

338 **Demo: Effect of MSOM Hierarchy on p**

339 The above plot is interesting because it shows that exceptionally high or exceptionally low chances to be
340 observed only occur for the species that were observed at least once; i.e., this says that if you didn't observe
341 it, it just takes on the mean. That's reasonable, I guess; but I also would think that the things that were
342 never observed could also be things that had a low chance of observability; but they could also have just a
343 low chance of actually being present. So I suppose in the end it just doesn't get informed, and reverts to the
344 mean?

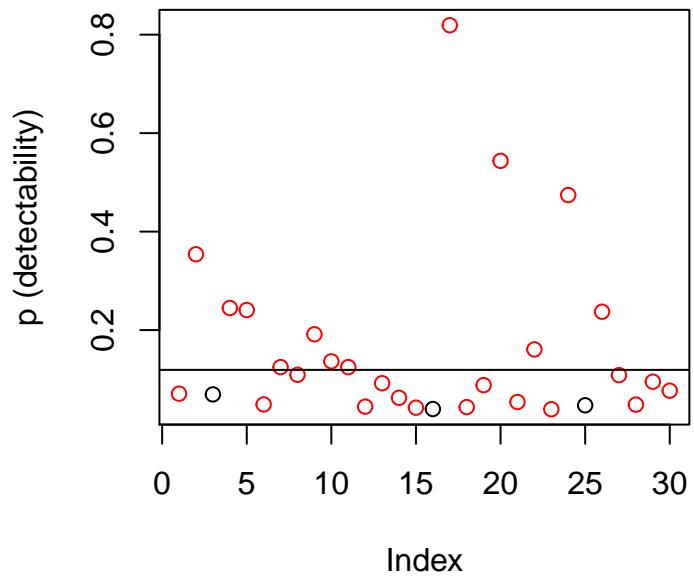


Figure 10: Probability of being detected, p . Horizontal line is mean probability. Figure only shows results for the first year of the simulation/ observation, and only 1 replicate. Different points are different species. Probability of being detected is a species-specific parameter (does not vary among sites, e.g.). Red points are species that were observed, black points are species that were never observed.

³⁴⁵ **Scatter Plot of \hat{p} vs p_{true}**

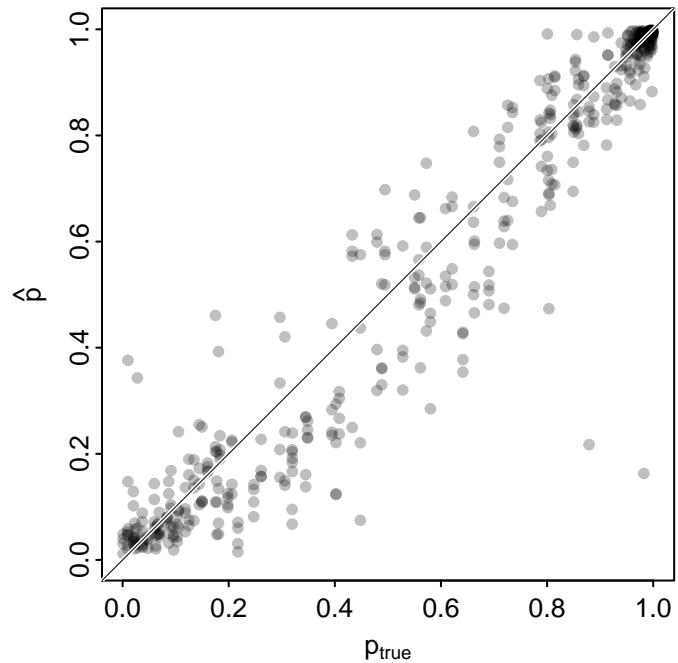


Figure 11: MSOM estimates (vertical axis) and true values of p_i , the species-specific (i) detection probability. Each point is subscripted by species i , year t , and observation replicate r .

³⁴⁶

³⁴⁷ **Scatter Plot of \hat{p} vs p_{true} , split by year and replicate**

³⁴⁸ Text explanation goes here

³⁴⁹

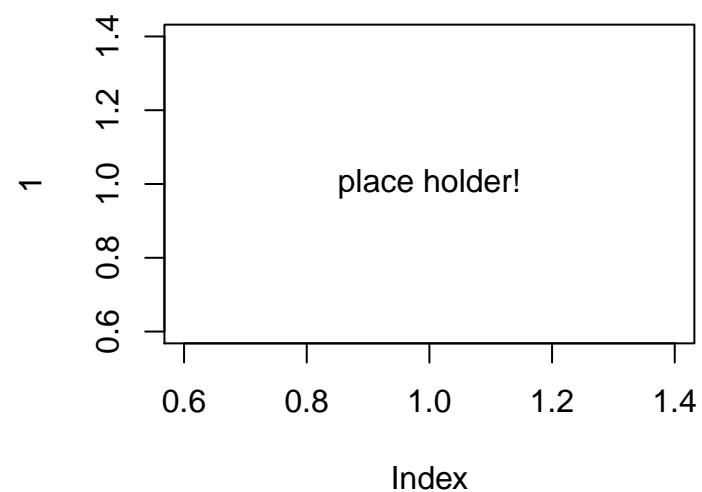


Figure 12: Caption goes here.

350 Assessment with Mixed Effects Models

351 Describe Motivation for Mixed Effects Models

- 352 **Motivation:** MSOM skill might differ across dimensions, trying to figure out what patterns I should expect
353 to pick out (spatial patterns in richness, temporal?) E.g., Is the correlation between MSOM and True the
354 same comparing across sites as comparing across years? Species, reps, also.
- 355 **Motivation:** What factors influence MSOM skill in a given dimension? E.g., Skill in finding differences in ψ
356 across species may depend on p , the chance of being identified. If p changes among years, might also explain
357 Read more about [specifying mixed effects models using lmer in R here](#)
- 358 This example is looking at ψ , probability of an individual species being present

359 Example LMER Analysis for ψ

```
# Just exploration/ starting point
library(car)
library(lme4)

blah <- reshape2:::melt.array(
  psi.true,
  varnames=c("site", "spp", "time", "rep"),
  value.name="true",
  as.is=T
)

blah.hat <- reshape2:::melt.array(
  psi.hat,
  varnames=c("site", "spp", "time", "rep"),
  value.name="hat",
  as.is=T
)

blah <- cbind(blah, hat=blah.hat[, "hat"])

blah$site <- as.factor(blah$site)
blah$spp <- as.factor(blah$spp)
blah$time <- as.factor(blah$time)
blah$rep <- as.factor(blah$rep)
```

```

# =====
# = Do LMER Analysis =
# =====

(blah.mod <- lmer(hat~true+(1|spp)+(1|time), data=blah))

```

```

360 ## Linear mixed model fit by REML ['lmerMod']
361 ## Formula: hat ~ true + (1 | spp) + (1 | time)
362 ## Data: blah
363 ## REML criterion at convergence: -7303.128
364 ## Random effects:
365 ## Groups Name Std.Dev.
366 ## spp (Intercept) 0.04946
367 ## time (Intercept) 0.02326
368 ## Residual 0.21984
369 ## Number of obs: 38880, groups: spp, 30; time, 4
370 ## Fixed Effects:
371 ## (Intercept) true
372 ## 0.1201 0.8734

```

```
Anova(blah.mod)
```

```

373 ## Analysis of Deviance Table (Type II Wald chisquare tests)
374 ##
375 ## Response: hat
376 ## Chisq Df Pr(>Chisq)
377 ## true 34126 1 < 2.2e-16 ***
378 ## ---
379 ## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

# =====
# = Explain LMER =
# =====

```

381 **Report Generation Notes**

382 **R Session Information**

```
383 ## R version 3.2.0 (2015-04-16)
384 ## Platform: x86_64-apple-darwin13.4.0 (64-bit)
385 ## Running under: OS X 10.10.5 (Yosemite)
386 ##
387 ## locale:
388 ## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
389 ##
390 ## attached base packages:
391 ## [1] parallel   grid      stats      graphics   grDevices  utils      datasets
392 ## [8] methods    base
393 ##
394 ## other attached packages:
395 ## [1] lme4_1.1-9       Matrix_1.2-0      car_2.0-26      kfigr_1.2
396 ## [5] xtable_1.7-4     rmarkdown_0.7     knitr_1.11     doParallel_1.0.8
397 ## [9] iterators_1.0.7   foreach_1.4.2    R2jags_0.5-6   rjags_3-15
398 ## [13] coda_0.17-1     igraph_0.7.1     fields_8.2-1   maps_2.3-9
399 ## [17] spam_1.0-1      data.table_1.9.4 raster_2.3-40  sp_1.1-0
400 ## [21] rbLib_0.0.2     taxize_0.5.2
401 ##
402 ## loaded via a namespace (and not attached):
403 ## [1] reshape2_1.4.1    splines_3.2.0    lattice_0.20-31 htmltools_0.2.6
404 ## [5] yaml_2.1.13      mgcv_1.8-6      chron_2.3-45   XML_3.98-1.3
405 ## [9] nloptr_1.0.4     plyr_1.8.2      stringr_1.0.0  codetools_0.2-11
406 ## [13] evaluate_0.7.2   Taxonstand_1.7  SparseM_1.7   permute_0.8-4
407 ## [17] quantreg_5.11   pbkrtest_0.4-2  numbers_0.6-1  highr_0.5
408 ## [21] Rcpp_0.11.6     formatR_1.2     vegan_2.3-0   jsonlite_0.9.16
409 ## [25] abind_1.4-3     digest_0.6.8    stringi_0.5-5 ssh.utils_1.0
410 ## [29] tools_3.2.0     bitops_1.0-6    magrittr_1.5   RCurl_1.95-4.6
411 ## [33] bold_0.2.6      cluster_2.0.1   ape_3.3       MASS_7.3-40
412 ## [37] minqa_1.2.4    assertthat_0.1  reshape_0.8.5 httr_0.6.1
413 ## [41] boot_1.3-16     R2WinBUGS_2.1-21 nnet_7.3-9   nlme_3.1-120
```

414 **Date Document Last Compiled**

```
415 ## Last compiled on: 2015-08-25
```