

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	15
²⁵	Definition of ψ	15
²⁶	Scatter Plot of Aggregated ψ	15
²⁷	Scatter Plot of $\hat{\psi}$ vs ψ_{true} , split by year and replicate	16
²⁸	Occupancy Response Curves	16
²⁹	Probability of Detection, p	20
³⁰	Definition of p	20
³¹	Demo: Effect of MSOM Hierarchy on p	20
³²	Scatter Plot of \hat{p} vs p_{true}	22
³³	Scatter Plot of \hat{p} vs p_{true} , split by year and replicate	22
³⁴	Assessment with Mixed Effects Models	24
³⁵	Describe Motivation for Mixed Effects Models	24
³⁶	Example LMER Analysis for ψ	24

37	Conclusion	29
38	Discussion of Results	29
39	Next Steps	29
40	Concluding Remarks	29
41	Report Generation Notes	30
42	R Session Information	30
43	Date Document Last Compiled	30

44

⁴⁵ Introduction

⁴⁶ Overview

⁴⁷ As water temperatures change, species may shift the size and location of their geographical ranges, bearing
⁴⁸ consequences for the food webs and economies linked to those species. However, species don't always respond
⁴⁹ similarly to shifting temperatures (different thermal tolerances, e.g.), which means that changing temperature
⁵⁰ may remix the composition and diversity of ecological communities.

⁵¹ The biological, spatial, and temporal scale of community diversity shifting in response to climate is massive.
⁵² A functional definition of a community may consist of 100's or 1000's of species, each of which may be
⁵³ shifting its range at a scale of decades and 100's kilometers. As a result, we need statistical methods for
⁵⁴ estimating biodiversity that don't rely on heavy replication and that make efficient use of available data.
⁵⁵ Enter the superstars: on the data side the trawl data set has amazing spatiotemporal and taxonomic extent
⁵⁶ and resolution; on the statistical side multispecies occupancy models (MSOM) are hierarchical state space
⁵⁷ models that are designed to estimate species richness and don't require consistent or extensive "replication".
⁵⁸ Although they're superstars, even these data and models have their limitations and pitfalls.

⁵⁹ Can we estimate the dynamics of species richness from trawl data using an MSOM? It's a hard question to
⁶⁰ answer because we can never know the "truth" for sure, but we can get an idea of how reliable our analysis
⁶¹ is by simulating fake data, for which we know true values because we created them. The trawl data set is
⁶² generated by two distinct processes: Nature's data generating process (NDGP), and the process by which
⁶³ humans observe the result of NDGP. So we ask: to what extent is the accuracy of estimates from an MSOM
⁶⁴ dependent on characteristics of NDGP, and in particular, the way in which we observe the result of NDGP?
⁶⁵ The strategy for answering this question is to simulate fake data where we approximate Nature but gain
⁶⁶ knowledge of "truth", "observe" the results of the true process, then try to recover the true species richness
⁶⁷ from these simulated data.

⁶⁸ The Simulation

⁶⁹ The goal of this simulation was to use a very basic process to generate presences and absences of species in
⁷⁰ space and time. In this version of the simulation, there is no explicit connection between years (they are
⁷¹ independent). There is a modest spatial connection, because in the simulation an environmental variable
⁷² determines habitat suitability. I think of this environmental variable as temperature, and I filled a grid with
⁷³ temperatures that ranged from the coldest at the top of the grid (north) and the warmest at the bottom
⁷⁴ (south) and added random variation among columns in the same row (among longitudes at the same latitude).

⁷⁵ One level of the simulation mimics NDGP. In this level, NDGP is best characterized by ψ , which is the product
⁷⁶ of a temperature and species' response curves. I.e., temperatures were used to determine the suitability of
⁷⁷ each grid cell to each simulated species. This suitability is known as ψ throughout this document.

⁷⁸ A second level of the simulation mimics human observation of NDGP — what we do when we collect data.
⁷⁹ This process was simulated by assigning each species has a unique probability of being observed or "detected"
⁸⁰ (this variable is p). The observation process gets several attempts at observing a given species in a given grid
⁸¹ cell; think of this as subdividing each site into subsites, and when you visit each subsite you have probability

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

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

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

98 Multispecies Occupancy Models (MSOMs)

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

₁₂₀ varies among space because of the environmental variable, and in a given location it varies among species
₁₂₁ because of the regression parameters.

₁₂₂

123 Conventions and Settings

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

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

133 Dimension Conventions

134 Summary

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

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

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

- 141 • Sub-sites are only relevant to the “observation” process
- 142 • Each site has the same number of possible sub-sites, but the number of sub-sites observed can vary
- 143 • 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

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

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

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

156 5. Replicates ($r = 4$)

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

162 **In Code**

163 The MSOM analyzes each year_t–replicate_r combination independently. Parameters subscripted by these
 164 dimensions are derived from separate analyses.

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

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

176

177 **Settings**

178 **Simulation Settings**

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

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

```
185 ## Dimensions: 81, 30, 4
186 ## grid.h = 9
187 ## grid.w = 9
188 ## grid.t = 4
189 ##
190 ## Number Species Possible (ns):
191 ## 30
192 ## Total Species Richness:
```

```

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

```

207 In the MSOM, detectability (p_i) is determined in the form of a logistic regression, which currently only
208 has an intercept (v_0) as predictor (so just a mean). That intercept varies among species (i.e., $v_{0,i}$), and
209 that variation is generated by drawing each individual species's intercept ($v_{0,i}$) from a parent distribution:
210 $v_{0,i} \sim \mathcal{N}(\mu_{v_0}, \sigma_{v_0}^2)$. See [section about \$p\$](#) for more info.

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

211

212 Settings for JAGS & MSOM

nChains	nIter	n0s	nSamples
3	50000	10	500

213 In the table above, **nChains**, **nIter**, and **nSamples** are all variables that are strictly pertinent to the Bayesian
214 analysis carried out in JAGS. The **n0s** value refers to the the degree of “data augmentation”. In this process,
215 you add extra species to the data set, and say that they were never observed. For our purposes, this is
216 employed for purely technical reasons, although it can be used to extra further inferences about species
217 richness.
218 The posterior samples from JAGS consist of 500 samples (above). However, to save memory and hard drive
219 space, I have often only saved measures of central tendency for each of these. In this assessment, I have
220 performed all calculations on the **centralT=median** of the posterior samples.

222 **Species Richness**

223 **Definition of species richness**

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

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

234 None of these estimates of richness include the 10 species that were part of the “data augmented”/ “adding
235 0's” process. Richness values can either be true (true simulated NDGP; R^{true}), observed (true simulated
236 human observation of NDGP; R^{obs}), or MSOM estimates of one of those two (\hat{R}^{true} or \hat{R}^{obs}).

237

238 **Regional Richness**

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

241 **Richness Boxplots**

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

247

248 **Richness Time Series**

249 Text explanation goes here

250 Need explanations for how each panel was calculated.

251 1. R^{true} is straightforward

252

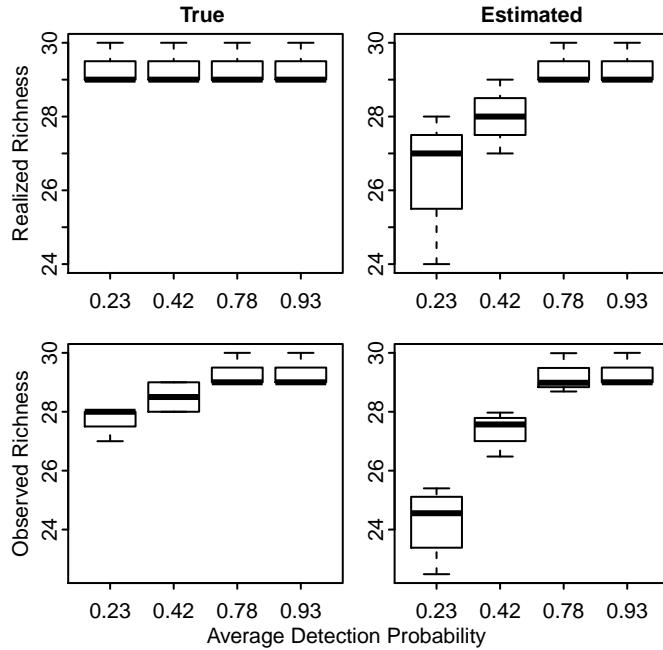


Figure 1: Boxplots of species richness. Numeric groupings indicate the average value of p across species during a given year-replicate combination. The panels in the left column are the true simulated values, and the panels on the right are the corresponding MSOM estimates. The top row indicates the 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.

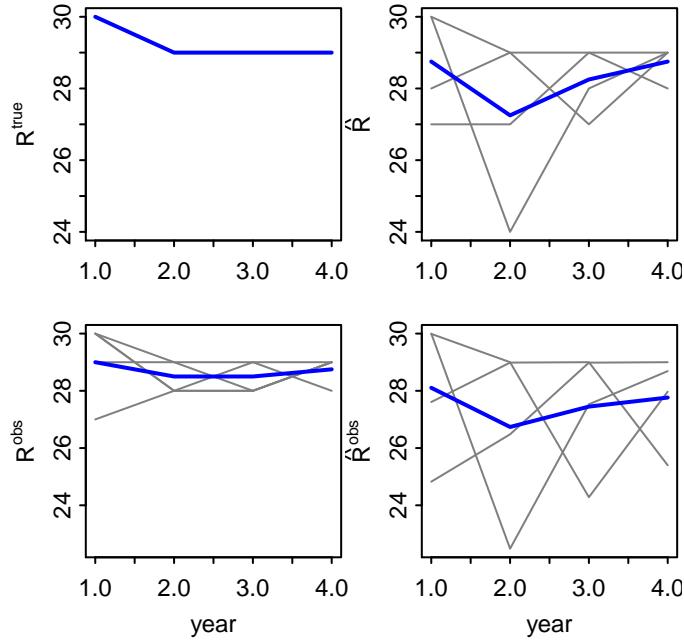


Figure 2: Time series of richness. Gray lines are individual replicates. Blue lines are averages. Note that detection probabilities ($p_{i,t,r}$, see [simulation settings above](#), as well as [definition of \$p\$ below](#)) change over time, and their temporal ordering differs among replicates.

253 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
 254 to the MSOM occurrence matrix (Y)

255

256 **Site Specific Richness (Nsite)**

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

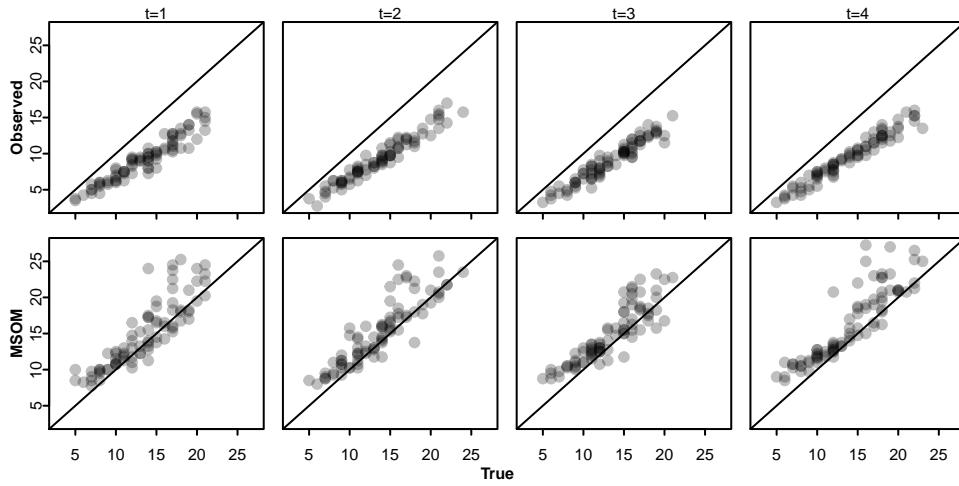


Figure 3: Site-specific richness (`Nsite`, N_j) from simulated observations (vertical axis, top row; N_j^{obs}) and from MSOM estimates (vertical axis, bottom row, \hat{N}_j) vs true site-specific richness (horizontal axis; N_j^*). The panel columns delineate the years of the simulation. Each point is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

258

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

260 Text explanation goes here

261

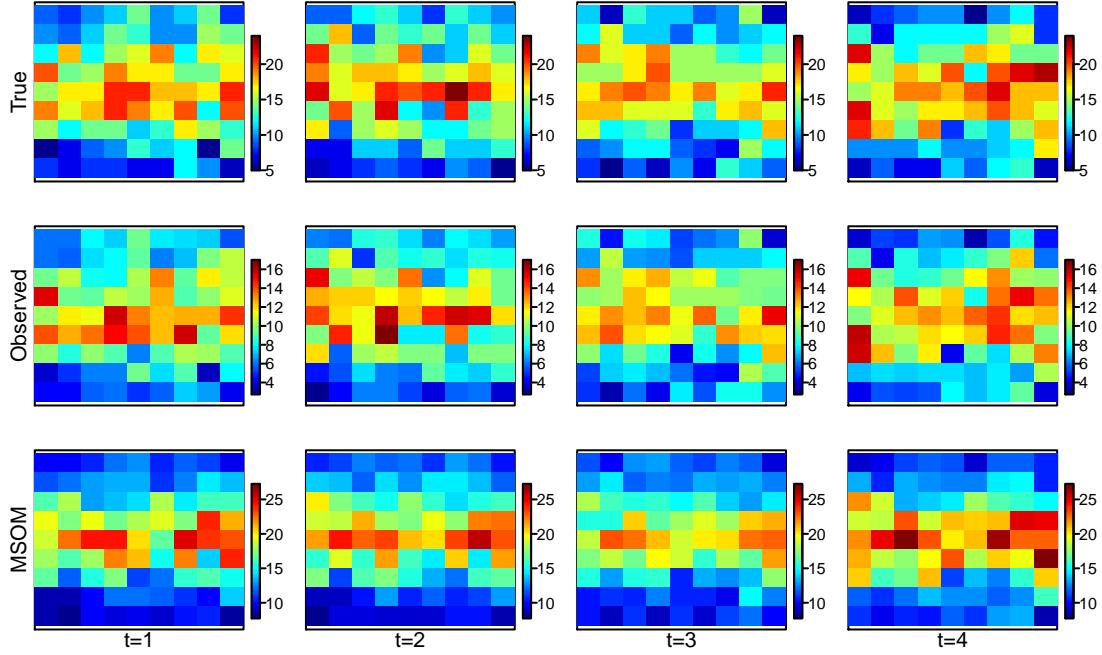


Figure 4: Maps of site- and year-specific species richness (N_{site}) from the simulation of the True process (top row), simulation of the Observed process (middle row), and the MSOM estimates (bottom row). X-axis and Y-axis indicate position in 2 dimensional space; it is important to note that the environmental 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.

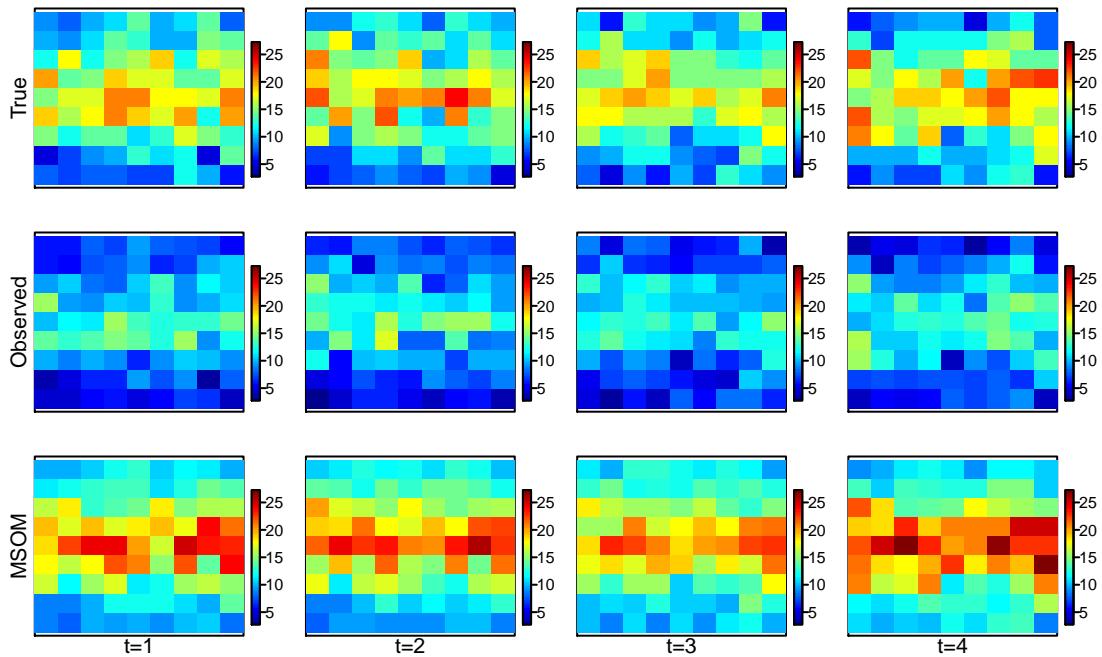


Figure 5: Same as previous figure, but all panels are on the same scale.

262 **Occupancy Probability, ψ**

263 **Definition of ψ**

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

268

269 **Scatter Plot of Aggregated ψ**

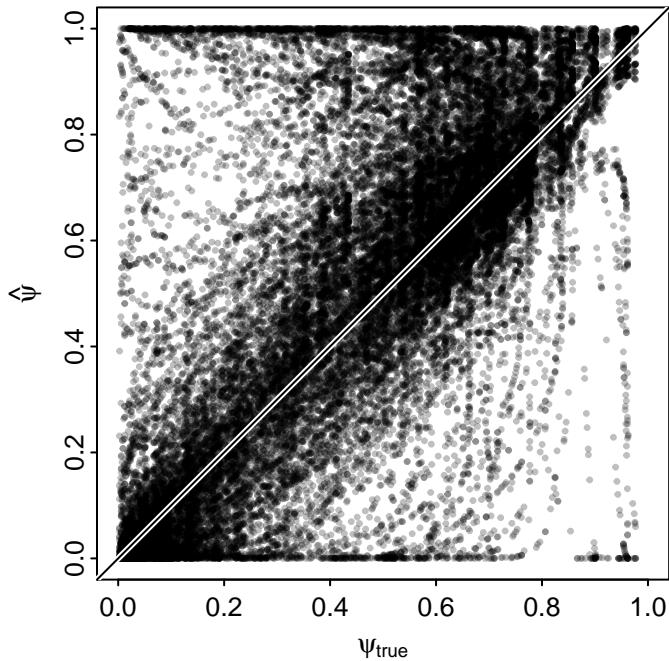


Figure 6: MSOM estimates of ψ ($\hat{\psi}$) vs. true values of ψ (ψ_{true}). Each point is a ψ value for a particular site-species-year-replicate. The white and black line is the 1:1 line.

- 270 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

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

274

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

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

277

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

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

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

291

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

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

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

298

299 Occupancy Response Curves

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

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

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

303 True Occupancy Response Curves

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

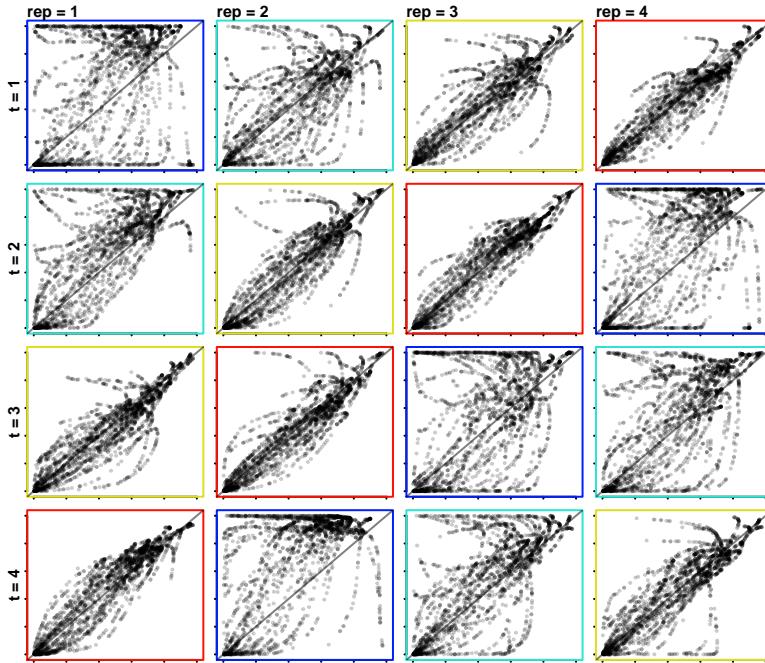


Figure 7: True (horizontal axes) and MSOM estimates (vertical axes) of occupancy probabilities ($\psi_{j,i,t,r}$) of species i occupying a location j . Years (t) are separated by rows, replicates (r) are separated by columns. The border color of each panel indicates the community-level mean probability 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.

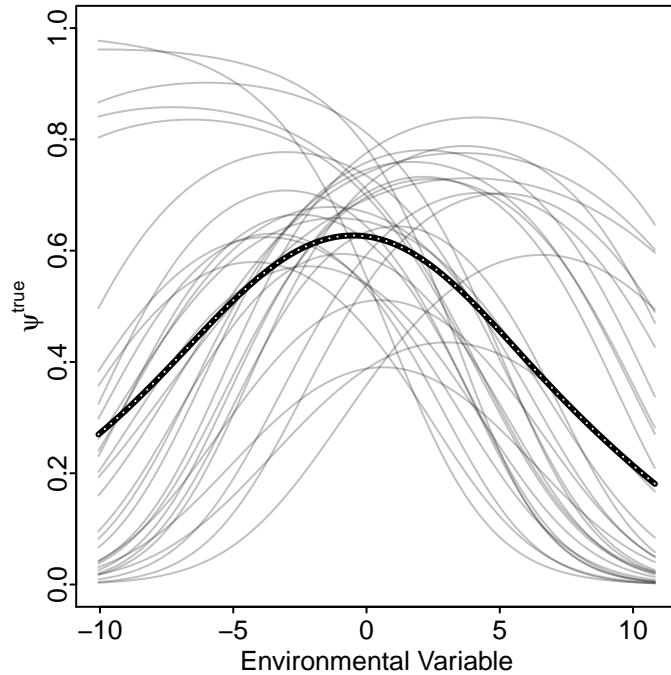


Figure 8: True simulated response curves. Vertical axis is the value of ψ^{true} , horizontal axis is the value of the environmental variable that, along with species-specific regression parameters, determines ψ^{true} . The thick line is the among-species mean value of ψ^{true} at a given value of the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

307 **Estimated Occupancy Response Curves**

308 ($\min(X) = -10.1$, and $\max(X) = 10.8$)

309 (red; $p_{max} = 0.93$)

310 (blue; $p_{min} = 0.23$)

311

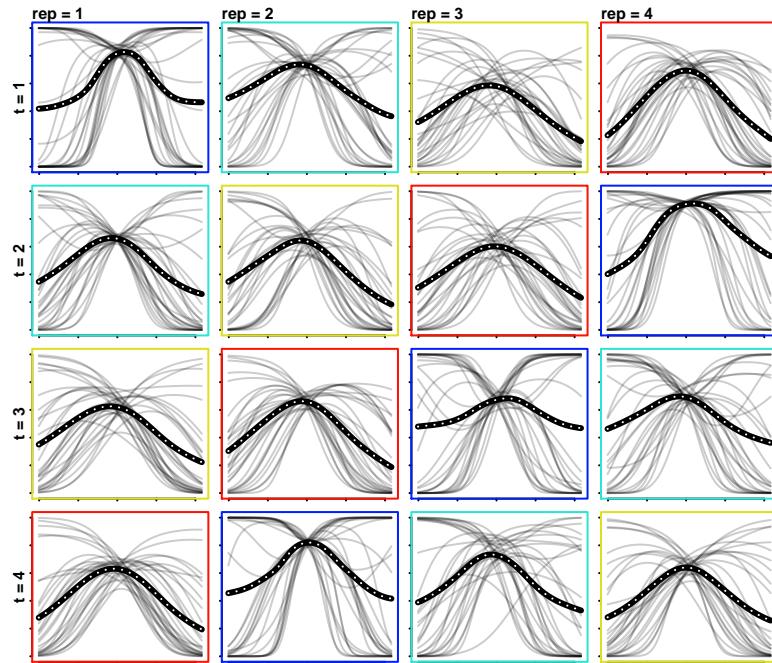


Figure 9: Response curves of species' probability of occupancy (ψ_i , vertical axis) across the full range of temperatures in the simulation. The color of the boxes around each panel refer to the among-species average of the probability of detection; warm colors indicate that the mean detection probability 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 the columns.

312 **Probability of Detection, p**

313 **Definition of p**

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

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

339

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

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

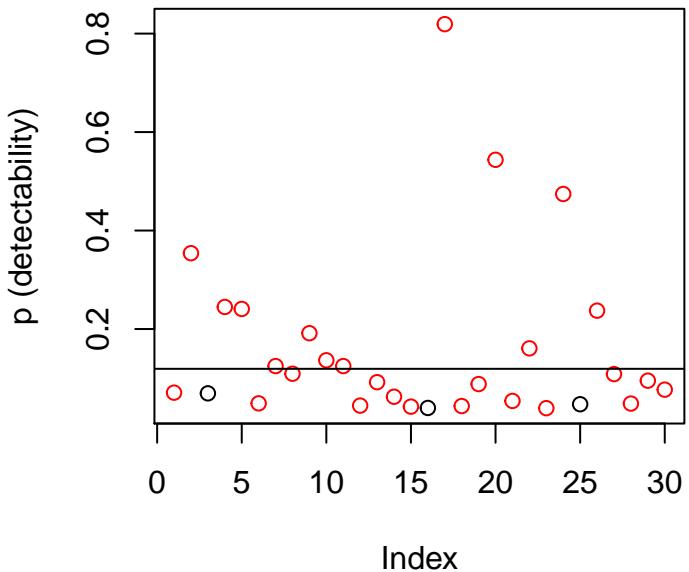


Figure 10: Probability of being detected, p . Horizontal line is mean probability. Figure only shows results for the first year of the simulation/ observation, and only 1 replicate. Different points are different species. Probability of being detected is a species-specific parameter (does not vary among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

³⁴⁷ **Scatter Plot of \hat{p} vs p_{true}**

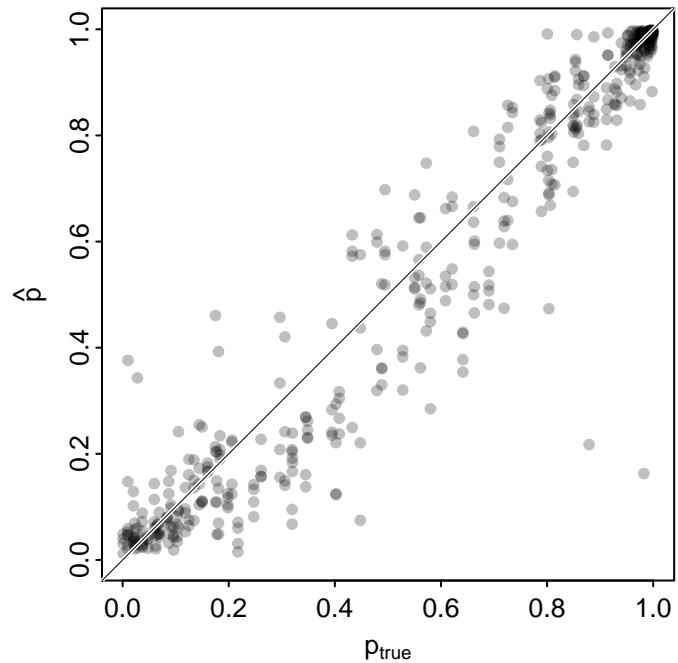


Figure 11: MSOM estimates (vertical axis) and true values of p_i , the species-specific (i) detection probability. Each point is subscripted by species i , year t , and observation replicate r .

³⁴⁸

³⁴⁹ **Scatter Plot of \hat{p} vs p_{true} , split by year and replicate**

³⁵⁰ Text explanation goes here

³⁵¹

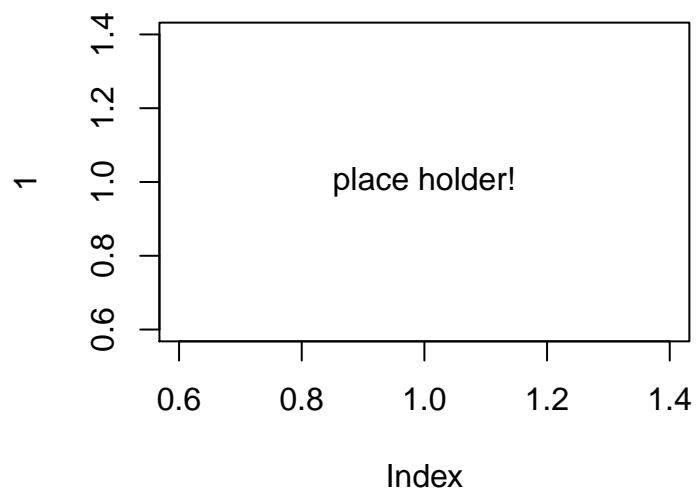


Figure 12: Caption goes here.

352 Assessment with Mixed Effects Models

353 Describe Motivation for Mixed Effects Models

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

361 Example LMER Analysis for ψ

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

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

# p
# p true
dat.p.true <- reshape2:::melt.array(
  aperm(array(p.true, dim=c(dim(p.true), dim(psi.true)[1])), c(4,1,2,3)),
  varnames=c("site", "spp", "time", "rep"),
  value.name="p.true",
```

```

    as.is=T
)
# p hat
dat.p.hat <- reshape2:::melt.array(
  aperm(array(p.hat, dim=c(dim(p.hat), dim(psi.hat)[1])), c(4,1,2,3)),
  varnames=c("site","spp","time","rep"),
  value.name="p.hat",
  as.is=T
)

# n.hauls
n.hauls <- sapply(big.out.obs, function(x)attributes(x)$n.haul)
n.hauls.dim <- c(grid.w*grid.h, n.obs.reps, ns, grid.t)
dat.n.hauls <- reshape2:::melt.array(
  aperm(array(n.hauls, dim=n.hauls.dim), c(1,3,4,2)),
  varnames=c("site","spp","time","rep"),
  value.name="n.hauls",
  as.is=T
)

# grid.X
# same structure (dims) as n.hauls
temp <- values(grid.X)
temp.dim <- c(grid.w*grid.h, n.obs.reps, ns, grid.t)
dat.temp <- reshape2:::melt.array(
  aperm(array(temp, dim=temp.dim), c(1,3,4,2)),
  varnames=c("site","spp","time","rep"),
  value.name="temp",
  as.is=T
)

# tax chance
tax.chance <- simplify2array(
  lapply(big.out.obs, function(x)(attributes(x)$obs.params)$tax.chance)
)
tax.chance.dim <- c(grid.t, ns, n.obs.reps, grid.w*grid.h)
dat.tax.chance <- reshape2:::melt.array(
  aperm(array(tax.chance, dim=tax.chance.dim), c(4,2,1,3)),
  varnames=c("site","spp","time","rep"),
  value.name="tax.chance",
  as.is=T
)

```

```

mod.dat <- cbind(
  dat.psi.true,
  psi.hat=dat.psi.hat[, "psi.hat"],
  p.true=dat.p.true[, "p.true"],
  p.hat=dat.p.hat[, "p.hat"],
  n.hauls=dat.n.hauls[, "n.hauls"],
  tax.chance=dat.tax.chance[, "tax.chance"],
  temp=dat.temp[, "temp"]
)
mod.dat[, "psi.error"] <- mod.dat[, "psi.hat"]-mod.dat[, "psi.true"]
mod.dat[, "p.error"] <- mod.dat[, "p.hat"] - mod.dat[, "p.true"]

mod.dat$site <- as.factor(mod.dat$site)
mod.dat$spp <- as.factor(mod.dat$spp)
mod.dat$time <- as.factor(mod.dat$time)
mod.dat$rep <- as.factor(mod.dat$rep)

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

mod1 <- lmer(psi.error~temp+(1|spp)+(1|site), data=mod.dat)
mod2 <- lmer(psi.error~n.hauls+(1|spp)+(1|site), data=mod.dat)
mod3 <- lmer(psi.error~p.error+(1|spp)+(1|site), data=mod.dat)
# mod4 <- lmer(psi.error~n.hauls-1+(1|spp)+(n.hauls-1|spp)+(1|site), data=mod.dat)
mod4 <- lmer(psi.error~temp+(1|spp)+(temp-1|spp)+(1|site), data=mod.dat)
mod5 <- lmer(psi.error~p.error+(1|spp)+(p.error-1|spp)+(1|site), data=mod.dat)
mod6 <- lmer(psi.error~p.error+(p.error|spp)+(1|site), data=mod.dat)

# Calculate covariance matrix (for spp)
mod6.varcor.spp <- attr(summary(mod6)$varcor$spp, "correlation")
mod6.varcor.spp <- format(mod6.varcor.spp, digits=2)
mod6.varcor.spp[!lower.tri(mod6.varcor.spp)] <- ""

mod.cap <- c(
  "Mixed effect models assessing sensitivity of $\psi_{\epsilon}$ to simulation conditions"
)

```

362 The goal with the mixed effects models was to understand what causes errors in ψ . I focused on ψ because it
 363 has all the information needed to understand variability in richness, but it has more information than the
 364 actual richness (richness is a community level statistic, ψ is species-specific). In these models, the response
 365 variable is $\hat{\psi} - \psi^{true}$, which we'll call ψ_ϵ . If we understand the source of variability in ψ_ϵ , then we can
 366 understand what leads to inaccuracies in our model.

Table 3: Mixed effect models assessing sensitivity of ψ_ϵ to simulation conditions

	Model 1	Model 2	Model 3	Model 4	Model 5	Model 6
(Intercept)	0.062*** (0.010)	0.060*** (0.011)	0.046*** (0.008)	0.062*** (0.010)	0.040*** (0.009)	0.040*** (0.009)
temp	0.001*** (0.000)			0.001 (0.002)		
n.hauls		0.001 (0.001)				
p.error			-0.618*** (0.012)		-0.867*** (0.136)	-0.868*** (0.136)
AIC	-6331.371	-6314.550	-8840.187	-8685.580	-11637.684	-11636.548
BIC	-6288.530	-6271.709	-8797.346	-8634.170	-11586.274	-11576.570
Log Likelihood	3170.686	3162.275	4425.094	4348.790	5824.842	5825.274
Num. obs.	38880	38880	38880	38880	38880	38880
Num. groups: site	81	81	81	81	81	81
Num. groups: spp	30	30	30	30	30	30
Variance: site.(Intercept)	0.000	0.000	0.000	0.000	0.000	0.000
Variance: spp.(Intercept)	0.003	0.003	0.002			0.002
Variance: Residual	0.049	0.049	0.046	0.046	0.043	0.043
Variance: spp.temp				0.000		
Variance: spp.1.(Intercept)				0.003	0.002	
Variance: spp.p.error					0.543	0.543

*** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$

367 When analyzing the trawl data, we will not know ψ_ϵ – we can obtain model residuals, but these are distinct
 368 from ψ_ϵ , because calculating ψ_ϵ requires knowing ψ^{true} which is a latent, unobserved variable. An empirical
 369 analysis would also lack some of the explanatory variables made available to us in the simulation. However, if
 370 we can explain variability in ψ_ϵ using simulated information that will also be available to the trawl analysis,
 371 then we can build intuition about the sources of error in our estimate of species richness even when we don't
 372 know the true value. And that is the fundamental goal of this document.

373 I'll highlight some of the things I learned from this analysis:

374 1. Neither the environmental variable (`temp`) nor the number of subsites sampled per site (`n.hauls`) were
 375 strongly related to ψ_ϵ

376 i. But `temp` might be a better predictor if transformed into an absolute value
 377 ii. When the model does not contain a site-specific random intercept (not shown here), `n.hauls`
 378 accounts for more variability

379

380 2. The `spp` random intercept explains much more variability than `site` equivalent

381

382 3. `p.error` explains a ton of variability

383 i. inversely related to ψ_ϵ ; intuition: if you have perfect detectability but you think it's terrible, you'll
 384 overestimate the true value

385 ii. model 6 is a worse fit than model 5, meaning that adding covariance structure between spp-specific
386 intercept and spp-specific p.error does not explain much variability

387 As stated above, a lot of variance is explained by the model term (p.error|spp), which allows the parameter
388 associated with p.error and the intercept (both fixed effects) to vary randomly among species. The
389 interpretation of these model terms is that

- 390 • Each species gets to draw its own intercept (“Variance: spp.1.(Intercept)” in the table) from a parent
391 distribution of intercepts
392 • A unit of error in the estimate of p has an influence on ψ_e that, similar to the intercept, varies among
393 species (“Variance: spp.p.error” in the table).

394 Therefore, the effect that a bad estimate of p has on ψ_e is not the same among species. It is not clear what
395 causes some species to be more sensitive to a poorly estimated p than others; one possibility is that p is
396 poorly estimated for species that have not been observed much, and it is this lack of observation that is also
397 responsible for generating uncertainty in $\hat{\psi}$. Regardless, a bad (good) estimate of p is a good predictor of a
398 bad (good) estimate of $\hat{\psi}$.

399

400 **Conclusion**

401 The MSOM reliably recover model parameters so long as average values of p_i were not at their lowest levels.

402 **Discussion of Results**

403 **Next Steps**

404 **Concluding Remarks**

405

406 **Report Generation Notes**

407 **R Session Information**

```
408 ## R version 3.1.2 (2014-10-31)
409 ## Platform: x86_64-apple-darwin13.4.0 (64-bit)
410 ##
411 ## locale:
412 ## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
413 ##
414 ## attached base packages:
415 ## [1] parallel   grid      stats      graphics   grDevices  utils      datasets
416 ## [8] methods    base
417 ##
418 ## other attached packages:
419 ## [1] texreg_1.35      memisc_0.97      MASS_7.3-35      lme4_1.1-6
420 ## [5] Rcpp_0.11.6       Matrix_1.1-4       car_2.0-24       rbLib_0.0.2
421 ## [9] stargazer_5.2     kfigr_1.2        xtable_1.7-4     rmarkdown_0.7
422 ## [13] knitr_1.10.5      doParallel_1.0.8  iterators_1.0.7 foreach_1.4.2
423 ## [17] R2jags_0.5-6      rjags_3-15       coda_0.16-1      lattice_0.20-29
424 ## [21] igraph_0.7.1      fields_6.9.1      maps_2.3-6       spam_0.41-0
425 ## [25] data.table_1.9.4  raster_2.3-24    sp_1.0-17
426 ##
427 ## loaded via a namespace (and not attached):
428 ## [1] abind_1.4-0        boot_1.3-17       chron_2.3-45
429 ## [4] codetools_0.2-9     digest_0.6.8       evaluate_0.7
430 ## [7] formatR_1.2         highr_0.5        htmltools_0.2.6
431 ## [10] mgcv_1.8-3         minqa_1.2.3      nlme_3.1-118
432 ## [13] nnet_7.3-8         numbers_0.5-6    pbkrtest_0.4-2
433 ## [16] plyr_1.8.1         quantreg_5.11    R2WinBUGS_2.1-19
434 ## [19] RcppEigen_0.3.2.1.1 reshape2_1.4.1  SparseM_1.6
435 ## [22] splines_3.1.2      ssh.utils_1.0    stringr_0.6.2
436 ## [25] tools_3.1.2        yaml_2.1.13
```

437 **Date Document Last Compiled**

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