

Multimodel approaches are not the best way to understand multifactorial systems

Benjamin M. Bolker*

25 July 2023

Abstract

Information-theoretic (IT) and multi-model averaging (MMA) statistical approaches are widely used but suboptimal tools for pursuing a multifactorial approach (also known as the method of multiple working hypotheses) in ecology. (1) Conceptually, IT encourages ecologists to perform tests on sets of artificial models. (2) MMA improves on IT model selection by implementing a simple form of *shrinkage estimation* (a way to make accurate predictions from a model with many parameters, by “shrinking” parameter estimates toward zero). However, other shrinkage estimators such as penalized regression or Bayesian hierarchical models with regularizing priors are more computationally efficient and better supported theoretically. (3) In general the procedures for extracting confidence intervals from MMA are overconfident, giving overly narrow intervals. If researchers want to accurately estimate the strength of multiple competing ecological processes along with reliable confidence intervals, the current best approach is to use full (maximal) statistical models after making principled, *a priori* decisions about which predictors to include.

Many modern ecological and evolutionary studies try to quantify the importance of multiple processes in ecological systems: for example, the effects of herbivory and fertilization on standing biomass (Gruner et al. 2008);

*Departments of Biology and Mathematics & Statistics, McMaster University, Hamilton, Ontario, Canada; bolker@mcmaster.ca; ORCID 0000-0002-2127-0443

26 the effects of bark, wood density, and fire on tree mortality (Brando et al.
27 2012); or the effects of taxonomic and genomic position on evolutionary rates
28 (Ghenu et al. 2016). This *multifactorial* approach (McGill 2016) complements,
29 rather than replacing, the traditional hypothesis-testing or strong-inferential
30 framework (Platt 1964; Fox 2016).¹

31 A standard approach to analyzing multifactorial systems, particularly com-
32 mon in wildlife and conservation ecology, goes as follows: (1) Construct a
33 full model that encompasses as many of the processes (and their interac-
34 tions) as is feasible. (2) Fit the full model and make sure that it describes
35 the data reasonably well (e.g. by computing R^2 values or estimating degree
36 of overdispersion). (3) Construct possible submodels of the full model by
37 setting subsets of parameters to zero. (4) Compute information-theoretic (IT)
38 measures of quality, such as the Akaike or Bayesian/Schwarz information
39 criteria (IC), for every submodel. (5) Use multi-model averaging (MMA) to
40 estimate model-averaged parameters and confidence intervals (CIs); possibly
41 draw conclusions about the importance of different processes by summing
42 the IC weights (Burnham and Anderson 2002). We argue that this approach,
43 even if used sensibly as advised by proponents of the approach (e.g. with
44 reasonable numbers of candidate submodels), is a poor way to approach
45 multifactorial problems.

46 Our goal is to tease apart the contributions of many processes, *all* of which
47 we believe are affecting our study system to some degree. If our scientific
48 question is (something like) “How important is this factor, in an absolute
49 sense or relative to other factors?”, not “Which of these factors are actually
50 doing *anything at all* in my system?”, why are we working so hard to fit
51 many models of which only one (the full model) incorporates all of the
52 factors? If we do not have particular, *a priori* discrete hypotheses (such
53 as “A influences the outcome but B does not”) about our system (and a
54 multifactorial approach would suggest that we should not), why does so
55 much of our data-analytic effort go into various ways to test between, or
56 combine and reconcile, multiple discrete models? In software engineering,
57 this would be called an “XY problem”²: rather than thinking about the best
58 way to solve our real problem *X* (understanding multifactorial systems), we
59 have gotten bogged down in the details of how to make a particular tool,
60 *Y* (multimodel approaches) provide the answers we need. Most critiques

¹While there is much interesting debate over the best methods for gathering evidence to distinguish among two or more particular, *intrinsically* discrete hypotheses (Taper and Ponciano 2015), that is not my focus here.

²<http://www.perlmonks.org/?node=XY+Problem>

61 of MMA address technical concerns such as the influence of unobserved
62 heterogeneity (Brewer, Butler, and Cooksley 2016), or criticize the misuse of
63 IT methods by ecologists (Cade 2015), but do not ask why we are comparing
64 discrete models in the first place. (Ecological statisticians are also beginning
65 to emphasize the importance of *causal inference* (Fieberg and Johnson 2015;
66 Laubach et al. 2021; Kimmel et al. 2021; Arif and MacNeil 2022); while this is
67 important, it is not the focus here.)

68 One legitimate reason to fit multiple models is as a step in a null-hypothesis
69 significance testing (NHST) procedure. While much maligned, NHSTs are a
70 useful part of data analysis — *not* to decide whether we really think a null
71 hypothesis is false (they almost always are), but to see if we can distinguish
72 signal from noise. Another interpretation is that NHSTs can test whether we
73 can reliably determine the *direction of effects* — that is, not whether the effect
74 of a predictor on some process is zero, but whether we can tell unequivocally
75 that it is positive (or negative, Jones and Tukey 2000). We can perform these
76 tests by statistically comparing a full model to a reduced model that pretends
77 the effect is exactly zero.

78 However, ecologists pursuing multimodel approaches are not fitting one-
79 step-reduced models to test hypotheses; they are fitting a wide range of
80 submodels, typically in the hope that multimodel averaging will help them
81 deal with insufficient data in a multifactorial world. If we had enough
82 information (even “big data” doesn’t always provide as much information
83 as we need), we could fit just the full model, drawing our conclusions from
84 the estimates and CIs with all of the factors considered simultaneously. But
85 we nearly always have too many predictors, and not enough data; we don’t
86 want to overfit (which will inflate our CIs and p-values to the point where we
87 can’t tell anything for sure), but at the same time we are scared of neglecting
88 potentially important effects.

89 Stepwise regression, the original strategy for separating signals from
90 noise, is now widely deprecated (Harrell 2001; Whittingham et al. 2006)³.
91 Information-theoretic tools mitigate the instability of stepwise approaches,
92 allow simultaneous comparison of many, non-nested models, and avoid the
93 stigma of NHST. A further step forward, multi-model averaging (Burnham
94 and Anderson 2002), accounts for model uncertainty and avoids focusing
95 on a single best model. Some forms of model averaging provide simple
96 *shrinkage estimators*; averaging the strength of effects between models where

³Although it may sometimes be adequate for selecting a single best model for prediction (Murtaugh 2009).

97 they are included and models where they are absent “shrinks” the estimated
98 effects toward zero (Cade 2015). More recently, however, model averaging
99 is experiencing a backlash, as statistically savvy ecologists point out that
100 multimodel averaging may run into trouble when variables are collinear
101 (Freckleton (2011; but cf. Walker 2017)); when we are careless about the
102 meaning of main effects in the presence of interactions; when we average
103 model parameters rather than model predictions (Cade 2015); or when we
104 use summed model weights to assess the relative importance of predictors
105 (Galipaud et al. (2014; but cf. Zhang, Zou, and Carroll 2015)).

106 IC were introduced into ecological science from applied ecology, by quanti-
107 tative ecologists who were focused on making the best possible predictions
108 to inform conservation and management. Now, however, rather than us-
109 ing IC as tools to identify the best predictive model, or to obtain the best
110 overall (model-averaged) predictions, most users of information criteria use
111 them either to quantify variable importance, or, by multimodel averaging, to
112 have their cake and eat it too — to avoid either over- or underfitting while
113 quantifying effects in multifactorial systems. Such users of IC encounter two
114 problems, one conceptual and one practical.

115 The conceptual problem with model averaging reflects the original sin of
116 unnecessarily discretizing a continuous world. Suppose we want to un-
117 derstand the effects of temperature and precipitation on biodiversity. The
118 model-comparison or model-averaging approach would construct five mod-
119 els: a null model with no effects of either temperature or precipitation, two
120 single-factor models, an additive model, and a full model allowing for in-
121 teractions between temperature and precipitation. We would then fit all (or
122 many) of these models and then model-average their parameters. We might
123 be doing this in an effort to get good predictions, or to test our confidence
124 that we know the signs of particular effects (measured in the context of
125 whatever processes are included in the reduced and the full models), but
126 they are only means to an end, and we shouldn’t fool ourselves into thinking
127 that we are using the method of multiple working hypotheses. For example,
128 Chamberlin (1897, reprinted as Raup and Chamberlin (1995)) argued that
129 in teaching about the origin of the Great Lakes we should urge students “to
130 conceive of three or more great agencies [pre-glacial erosion, glacial erosion,
131 crust deformation] working successively or simultaneously, and to estimate
132 how much was accomplished by each of these agencies.” Chamberlin was
133 *not* suggesting that we test which individual mechanism or combination
134 of mechanisms fits the data best (in whatever sense), but instead that we
135 acknowledge that the world is multifactorial.

136 The technical problem with model averaging is its computational inefficiency.
137 Individual models can take minutes or hours to fit, and we may have to
138 fit dozens or scores of sub-models in the multi-model averaging process.
139 There are efficient tools available for fitting “right-sized” models that avoid
140 many of the technical problems of model averaging. Penalized methods
141 such as ridge and lasso regression (Dahlgren 2010) are well known outside
142 of ecology; in a Bayesian setting, informative priors centered at zero have
143 the same effect of *regularizing* — pushing weak effects toward zero and
144 controlling model complexity (more or less synonymous with the *shrinkage*
145 of estimates described above). Developed specifically for optimal (predictive)
146 fitting in models with many parameters, these models have well-understood
147 statistical properties; they avoid the pitfalls of model-averaging correlated
148 or nonlinear parameters; and, by avoiding the need to fit many sub-models
149 in the model-averaging processes, they are much faster.⁴

150 Here I am not tackling the issue of whether ‘truth’ is included in our model
151 set (it isn’t), and how this matters to our inference (Bernardo and Smith
152 1994; Barker and Link 2015). I am claiming the opposite, that our full model
153 is usually as close to truth as we can get; we don’t really believe any of
154 the less complex models. If we are trying to get the best predictions, or to
155 compare the strength of various processes in a multifactorial context, there
156 may be better ways to do it. In situations where we really want to compare
157 qualitatively different, non-nested hypotheses (Luttbeg, Langen, and Adams
158 2004), AIC or BIC or any appropriate model-comparison tool is fine; however,
159 if the models are *really* qualitatively different, perhaps we shouldn’t be trying
160 to merge them by averaging.

161 Penalized models have their own difficulties. A big advantage of IC-based
162 methods is that, like wrapper methods for feature selection in machine
163 learning (Chandrashekar and Sahin 2014), we can use model averaging as
164 long as we can fit component models and extract the log-likelihood and
165 number of parameters — we never need to develop any additional soft-
166 ware. Although powerful computational tools exist for fitting penalized
167 versions of linear and generalized linear models (e.g. the `glmnet` package
168 for R) and mixed models (`glmmLasso`), software for some of the more exotic
169 models used by ecologists (e.g. zero-inflated models) may not be readily
170 available. Fitting these models requires the user to choose the degree of
171 penalization. Although this process is conveniently automated in tools like

⁴Although they often require a computationally expensive cross-validation step in order to choose the degree of penalization.

172 `glmnet`, correctly assessing out-of-sample accuracy (and hence the correct
173 level of penalization) is tricky for data that are correlated in space or time
174 (Wenger and Olden 2012; Roberts et al. 2016).

175 Finally, inference (computing p-values and CIs) for parameters in penalized
176 models — one of the most basic outputs we need from a statistical analysis
177 of a multifactorial system — is a current research problem; statisticians have
178 proposed a variety of methods (Pötscher and Schneider 2010; Javanmard
179 and Montanari 2014; Lockhart et al. 2014; Taylor and Tibshirani 2018), but
180 they are far from being standard options in software. Ecologists should en-
181 courage their quantitatively savvy friends to build tools that make penalized
182 approaches easier to use.

183 Statisticians derived confidence intervals for ridge regression long ago (Oben-
184 chain 1977) — but, surprisingly, they are identical to the confidence intervals
185 one would have gotten from the full model without penalization! Wang and
186 Zhou (2013) similarly proved that model-averaging CIs derived as suggested
187 by Hjort and Claeskens (2003) are asymptotically (i.e. for arbitrarily large data
188 sets) equivalent to the CIs from the full model. Analytical and simulation
189 studies (Turek and Fletcher 2012; Fletcher and Turek 2012; Turek 2013, 2015;
190 Kabaila, Welsh, and Abeysekera 2016; Dormann et al. 2018) have shown
191 that a variety of alternative methods for constructing CIs are overoptimistic,
192 i.e. that they generate too-narrow confidence intervals with coverage lower
193 than the nominal level. Simulations from several of the studies above show
194 that MMA confidence intervals constructed according to the best known
195 procedures typically include the true parameter values only about 80% or
196 90% of the time. In particular, Kabaila, Welsh, and Abeysekera (2016) say that
197 constructing CIs that take advantage of shrinkage but still achieve correct
198 coverage will be very difficult to achieve using model averaged confidence
199 intervals. (The only examples I have been able to find of MMA confidence
200 intervals with close to nominal coverage are from Chapter 5 of Burnham
201 and Anderson (2002).) In short, it seems difficult to find model-averaged
202 confidence intervals that compete successfully with the standard confidence
203 interval based on the full model.

204 Free lunches do not exist in statistics, any more than anywhere else. We can
205 use penalized approaches to improve prediction accuracy without having
206 to sacrifice any input variables (by trading bias for variance), but the only
207 known way to gain statistical power for testing hypotheses, or narrowing
208 our uncertainty about our predictions, is to limit the scope of our models *a*
209 *priori* (Harrell 2001), to add information from pre-specified Bayesian priors

(or equivalent regularization procedures), or to collect more data.

If we have good experimental designs and sensible scientific questions, muddling through with existing techniques will often give reasonable results (Murtaugh 2009). But ecologists should at the very least be aware that the roundabout statistical methods they currently rely on to understand multifactorial systems were designed for prediction rather than inference. When prediction is the primary goal, penalized methods can work better (faster and with better-understood statistical properties) than multimodel averaging. When estimating the magnitude of effects or judging variable importance, penalized methods may be appropriate — or we may have to go back to the difficult choice of focusing on a restricted number of variables for which we have enough data to fit and interpreting the full model.

Acknowledgements

Thanks to Jonathan Dushoff for conversations on these topics over many years. Dana Karelus, Daniel Turek, and Jeff Walker provided useful input: Noam Ross encouraged me to finally submit the paper; Tara Bolker gave advice on straw men. This work was supported by multiple NSERC Discovery grants.

References

- Arif, Suchinta, and M. Aaron MacNeil. 2022. "Predictive Models Aren't for Causal Inference." *Ecology Letters* 25 (8): 1741–5. <https://doi.org/10.1111/ele.14033>.
- Barker, Richard J., and William A. Link. 2015. "Truth, Models, Model Sets, AIC, and Multimodel Inference: A Bayesian Perspective." *The Journal of Wildlife Management* 79 (5): 730–38. <https://doi.org/10.1002/jwmg.890>.
- Bernardo, José M., and Adrian F. M. Smith. 1994. *Bayesian Theory*. 1st ed. John Wiley & Sons, Ltd. <https://doi.org/10.1002/9780470316870>.
- Brando, P. M., D. C. Nepstad, J. K. Balch, B. Bolker, M. C. Christman, M. Coe, and F. E. Putz. 2012. "Fire-Induced Tree Mortality in a Neotropical Forest: The Roles of Bark Traits, Tree Size, Wood Density and Fire Behavior." *Global Change Biology* 18 (2): 630–41. <https://doi.org/10.1111/j.1365-2486.2011.02533.x>.

- 242 Brewer, Mark J., Adam Butler, and Susan L. Cooksley. 2016. "The Rela-
243 tive Performance of AIC, AICC and BIC in the Presence of Unobserved
244 Heterogeneity." *Methods in Ecology and Evolution* 7 (6): 679–92. <https://doi.org/10.1111/2041-210X.12541>.
245
- 246 Burnham, Kenneth P., and David R. Anderson. 2002. *Model Selection and*
247 *Multimodel Inference: A Practical Information-Theoretic Approach*. Springer.
- 248 Cade, Brian S. 2015. "Model Averaging and Muddled Multimodel Inference."
249 *Ecology*. <https://doi.org/10.1890/14-1639.1>.
- 250 Chandrashekar, Girish, and Ferat Sahin. 2014. "A Survey on Feature
251 Selection Methods." *Computers & Electrical Engineering* 40 (1): 16–28.
252 <https://doi.org/10.1016/j.compeleceng.2013.11.024>.
- 253 Dahlgren, Johan P. 2010. "Alternative Regression Methods Are Not Consid-
254 ered in Murtaugh (2009) or by Ecologists in General." *Ecology Letters* 13
255 (5): E7–E9. <https://doi.org/10.1111/j.1461-0248.2010.01460.x>.
- 256 Dormann, Carsten F., Justin M. Calabrese, Gurutzeta Guillera-Arroita, Eleni
257 Matechou, Volker Bahn, Kamil Bartoń, Colin M. Beale, et al. 2018. "Model
258 Averaging in Ecology: A Review of Bayesian, Information-Theoretic and
259 Tactical Approaches for Predictive Inference." *Ecological Monographs*.
260 <https://doi.org/10.1002/ecm.1309>.
- 261 Fieberg, John, and Douglas H. Johnson. 2015. "MMI: Multimodel Inference
262 or Models with Management Implications?: Multimodel Inference and
263 Models for Management." *The Journal of Wildlife Management* 79 (5): 708–
264 18. <https://doi.org/10.1002/jwmg.894>.
- 265 Fletcher, David, and Daniel Turek. 2012. "Model-Averaged Profile Like-
266 lihood Intervals." *Journal of Agricultural, Biological, and Environmental*
267 *Statistics* 17 (1): 38–51.
- 268 Fox, Jeremy. 2016. "Why Don't More Ecologists Use Strong Inference?"
269 *Dynamic Ecology*. <https://dynamicecology.wordpress.com/2016/06/01/obstacles-to-strong-inference-in-ecology/>.
270
- 271 Freckleton, Robert P. 2011. "Dealing with Collinearity in Behavioural and
272 Ecological Data: Model Averaging and the Problems of Measurement
273 Error." *Behavioral Ecology and Sociobiology* 65 (1): 91–101.
- 274 Galipaud, Matthias, Mark A. F. Gillingham, Morgan David, and François-
275 Xavier Dechaume-Moncharmont. 2014. "Ecologists Overestimate the

276 Importance of Predictor Variables in Model Averaging: A Plea for Cau-
 277 tious Interpretations." *Methods in Ecology and Evolution* 5 (10): 983–91.
 278 <https://doi.org/10.1111/2041-210X.12251>.

279 Ghenu, Ana-Hermina, Benjamin M. Bolker, Don J. Melnick, and Ben J. Evans.
 280 2016. "Multicopy Gene Family Evolution on Primate Y Chromosomes."
 281 *BMC Genomics* 17: 157. <https://doi.org/10.1186/s12864-015-2187-8>.

282 Gruner, D. S., J. E. Smith, E. W. Seabloom, S. A. Sandin, J. T. Ngai, H. Hille-
 283 brand, W. S. Harpole, et al. 2008. "A Cross-System Synthesis of Consumer
 284 and Nutrient Resource Control on Producer Biomass." *Ecology Letters* 11
 285 (7): 740–55.

286 Harrell, Frank. 2001. *Regression Modeling Strategies*. Springer.

287 Hjort, Nils Lid, and Gerda Claeskens. 2003. "Frequentist Model Average
 288 Estimators." *Journal of the American Statistical Association* 98 (464): 879–99.
 289 <https://doi.org/10.1198/016214503000000828>.

290 Javanmard, Adel, and Andrea Montanari. 2014. "Confidence Intervals and
 291 Hypothesis Testing for High-Dimensional Regression." *The Journal of*
 292 *Machine Learning Research* 15 (1): 2869–2909. <http://dl.acm.org/citation.cfm?id=2697057>.

294 Jones, Lyle V., and John W. Tukey. 2000. "A Sensible Formulation of the
 295 Significance Test." *Psychological Methods* 5 (4): 411–14. <https://doi.org/10.1037//1082-989X.5.4.411>.

297 Kabaila, Paul, A. H. Welsh, and Waruni Abeysekera. 2016. "Model-Averaged
 298 Confidence Intervals." *Scandinavian Journal of Statistics* 43 (1): 35–48.
 299 <https://doi.org/10.1111/sjos.12163>.

300 Kimmel, Kaitlin, Laura E. Dee, Meghan L. Avolio, and Paul J. Ferraro. 2021.
 301 "Causal Assumptions and Causal Inference in Ecological Experiments."
 302 *Trends in Ecology & Evolution* 36 (12): 1141–52. [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.tree.2021.08.008)
 303 [tree.2021.08.008](https://doi.org/10.1016/j.tree.2021.08.008).

304 Laubach, Zachary M., Eleanor J. Murray, Kim L. Hoke, Rebecca J. Safran,
 305 and Wei Perng. 2021. "A Biologist's Guide to Model Selection and Causal
 306 Inference." *Proceedings of the Royal Society B: Biological Sciences* 288 (1943):
 307 20202815. <https://doi.org/10.1098/rspb.2020.2815>.

308 Lockhart, Richard, Jonathan Taylor, Ryan J. Tibshirani, and Robert Tibshirani.
 309 2014. "A Significance Test for the Lasso." *Annals of Statistics* 42 (2): 413.
 310 <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4285373/>.

- 311 Luttbegg, Barney, Tom A. Langen, and Associate Editor: Eldridge S. Adams.
312 2004. "Comparing Alternative Models to Empirical Data: Cognitive
313 Models of Western Scrub-Jay Foraging Behavior." *The American Naturalist*
314 163 (2): 263–76. <https://doi.org/10.1086/381319>.
- 315 McGill, Brian. 2016. "Why Ecology Is Hard (and Fun) – Multicausality."
316 *Dynamic Ecology*. <https://dynamicecology.wordpress.com/2016/03/02/why-ecology-is-hard-and-fun-multicausality/>.
- 317
318 Murtaugh, Paul A. 2009. "Performance of Several Variable-Selection Meth-
319 ods Applied to Real Ecological Data." *Ecology Letters* 12 (10): 1061–8.
320 <https://doi.org/10.1111/j.1461-0248.2009.01361.x>.
- 321 Obenchain, R. 1977. "Classical F -Tests and Confidence Regions for Ridge
322 Regression." *Technometrics* 19 (4): 429–39.
- 323 Platt, John R. 1964. "Strong Inference." *Science*, New Series, 146 (3642):
324 347–53. <https://doi.org/10.2307/1714268>.
- 325 Pötscher, Benedikt M., and Ulrike Schneider. 2010. "Confidence Sets Based
326 on Penalized Maximum Likelihood Estimators in Gaussian Regression."
327 *Electronic Journal of Statistics* 4 (January): 334–60. <https://doi.org/10.1214/09-EJS523>.
328
- 329 Raup, David C., and T. C. Chamberlin. 1995. "The Method of Multiple
330 Working Hypotheses." *The Journal of Geology* 103 (3): 349–54. <http://www.jstor.org/stable/30071227>.
331
- 332 Roberts, David R., Volker Bahn, Simone Ciuti, Mark S. Boyce, Jane Elith,
333 Gurutzeta Guillera-Aroita, Severin Hauenstein, et al. 2016. "Cross-
334 Validation Strategies for Data with Temporal, Spatial, Hierarchical, or
335 Phylogenetic Structure." *Ecography*, December, 913–29. <https://doi.org/10.1111/ecog.02881>.
336
- 337 Taper, Mark L., and José Miguel Ponciano. 2015. "Evidential Statistics as a
338 Statistical Modern Synthesis to Support 21st Century Science." *Population*
339 *Ecology* 58 (1): 9–29. <https://doi.org/10.1007/s10144-015-0533-y>.
- 340 Taylor, Jonathan, and Robert Tibshirani. 2018. "Post-Selection Inference for
341 L1-penalized Likelihood Models." *Canadian Journal of Statistics* 46 (1):
342 41–61. <https://doi.org/10.1002/cjs.11313>.
- 343 Turek, Daniel. 2015. "Comparison of the Frequentist MATA Confidence
344 Interval with Bayesian Model-Averaged Confidence Intervals." *Journal of*
345 *Probability and Statistics* 2015. <https://doi.org/10.1155/2015/420483>.

- 346 Turek, Daniel Bernard. 2013. "Frequentist Model-Averaged Confidence
347 Intervals." PhD thesis, University of Otago. [https://www.otago.ourarc
348 hive.ac.nz/bitstream/handle/10523/3923/TurekDanielB2013PhD.pdf](https://www.otago.ourarchive.ac.nz/bitstream/handle/10523/3923/TurekDanielB2013PhD.pdf).
- 349 Turek, Daniel, and David Fletcher. 2012. "Model-Averaged Wald Confidence
350 Intervals." *Computational Statistics & Data Analysis* 56 (9): 2809–15. [https:
351 //doi.org/10.1016/j.csda.2012.03.002](https://doi.org/10.1016/j.csda.2012.03.002).
- 352 Walker, Jeffrey A. 2017. "A Defense of Model Averaging." *bioRxiv*, 133785.
353 <https://doi.org/10.1101/133785>.
- 354 Wang, Haiying, and Sherry Z. F. Zhou. 2013. "Interval Estimation by Frequen-
355 tist Model Averaging." *Communications in Statistics - Theory and Methods*
356 42 (23): 4342–56. <https://doi.org/10.1080/03610926.2011.647218>.
- 357 Wenger, Seth J., and Julian D. Olden. 2012. "Assessing Transferability of
358 Ecological Models: An Underappreciated Aspect of Statistical Validation."
359 *Methods in Ecology and Evolution* 3 (2): 260–67. [https://doi.org/10.1111/j.
360 2041-210X.2011.00170.x](https://doi.org/10.1111/j.2041-210X.2011.00170.x).
- 361 Whittingham, Mark J., Philip A. Stephens, Richard B. Bradbury, and Robert P.
362 Freckleton. 2006. "Why Do We Still Use Stepwise Modelling in Ecology
363 and Behaviour?" *Journal of Animal Ecology* 75 (5): 1182–9. [https://doi.or
364 g/10.1111/j.1365-2656.2006.01141.x](https://doi.org/10.1111/j.1365-2656.2006.01141.x).
- 365 Zhang, Xinyu, Guohua Zou, and Raymond J. Carroll. 2015. "Model Aver-
366 aging Based on Kullback-Leibler Distance." *Statistica Sinica* 25: 1583–98.
367 <https://doi.org/10.5705/ss.2013.326>.