

**A peer-reviewed version of this preprint was published in PeerJ on 23 May 2018.**

[View the peer-reviewed version](https://peerj.com/articles/4794) (peerj.com/articles/4794), which is the preferred citable publication unless you specifically need to cite this preprint.

Harrison XA, Donaldson L, Correa-Cano ME, Evans J, Fisher DN, Goodwin CED, Robinson BS, Hodgson DJ, Inger R. 2018. A brief introduction to mixed effects modelling and multi-model inference in ecology. PeerJ 6:e4794 <https://doi.org/10.7717/peerj.4794>

1 A Brief Introduction to Mixed Effects Modelling and Multi-model Inference in Ecology

2

3 Xavier A. Harrison<sup>1</sup>, Lynda Donaldson<sup>2</sup>, Maria Eugenia Correa-Cano<sup>2</sup>, Julian Evans<sup>3,4</sup>,  
4 David N. Fisher<sup>3&5</sup>, Cecily E. D. Goodwin<sup>2</sup>, Beth S. Robinson<sup>2&6</sup>, David J. Hodgson<sup>3</sup> and  
5 Richard Inger<sup>2&3</sup>.

6

7 <sup>1</sup> Institute of Zoology, Zoological Society of London, London, United Kingdom

8 <sup>2</sup> Environment and Sustainability Institute, University of Exeter, Penryn, United Kingdom

9 <sup>3</sup> Centre for Ecology and Conservation, University of Exeter, Penryn, United Kingdom

10 <sup>4</sup> Department of Biology, University of Ottawa, Ottawa, Canada

11 <sup>5</sup> Department of Integrative Biology, University of Guelph, Guelph, Canada

12 <sup>6</sup> WildTeam Conservation, Padstow, United Kingdom

13

14

15 Corresponding Authors:

16 Xavier Harrison xav.harrison@gmail.com

17 Richard Inger rich.inger@gmail.com

18

19

20

21

22

23

24

25

26

27

28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50

## ABSTRACT

The use of linear mixed effects models (LMMs) is increasingly common in the analysis of biological data. Whilst LMMs offer a flexible approach to modelling a broad range of data types, ecological data are often complex and require complex model structures, and the fitting and interpretation of such models is not always straightforward. The ability to achieve robust biological inference requires that practitioners know how and when to apply these tools. Here, we provide a general overview of current methods for the application of LMMs to biological data, and highlight the typical pitfalls that can be encountered in the statistical modelling process. We tackle several issues relating to the use of information theory and multi-model inference in ecology, and demonstrate the tendency for data dredging to lead to greatly inflated Type I error rate (false positives) and impaired inference. We offer practical solutions and direct the reader to key references that provide further technical detail for those seeking a deeper understanding. This overview should serve as a widely accessible code of best practice for applying LMMs to complex biological problems and model structures, and in doing so improve the robustness of conclusions drawn from studies investigating ecological and evolutionary questions.

51

## 52 Introduction

53

54 In recent years, the suite of statistical tools available to biologists and the complexity of  
55 biological data analyses have grown in tandem (Low-Decarie et al 2014; Zuur & Ieno  
56 2016; Kass et al 2016). The availability of novel and sophisticated statistical techniques  
57 means we are better equipped than ever to extract signal from noisy biological data, but  
58 it remains challenging to know how to apply these tools, and which statistical  
59 technique(s) might be best suited to answering specific questions (Kass et al 2016).  
60 Often, simple analyses will be sufficient (Murtaugh 2007), but more complex data  
61 structures often require more complex methods such as linear mixed effects models  
62 (Zuur et al 2009), generalized additive models (Wood et al 2015) or Bayesian inference  
63 (Ellison 2004). Both accurate parameter estimates and robust biological inference  
64 require that ecologists be aware of the pitfalls and assumptions that accompany these  
65 techniques and adjust modelling decisions accordingly (Bolker et al 2009).

66 Linear mixed effects models (LMMs) and generalized linear mixed effects models  
67 (GLMMs), have gained significant traction in the last decade (Zuur et al 2009; Bolker et  
68 al 2009). Both extend traditional linear models to include a combination of fixed and  
69 random effects as predictor variables. The introduction of random effects affords several  
70 non-exclusive benefits. First, biological datasets are often highly structured, containing  
71 clusters of non-independent observational units that are hierarchical in nature, and  
72 LMMs allow us to explicitly model the non-independence in such data. For example, we  
73 might measure several chicks from the same clutch, and several clutches from different  
74 females, or we might take repeated measurements of the same chick's growth rate over  
75 time. In both cases, we might expect that measurements within a statistical unit (here,  
76 an individual, or a female's clutch) might be more similar than measurements from  
77 different units. Explicit modelling of the random effects structure will aid correct  
78 inference of fixed effects, depending on which level of the system's hierarchy is being  
79 manipulated. In our example, if the fixed effect varies or is manipulated at the level of

3

80 the clutch, then pseudoreplicated measurements of each chick can be controlled  
81 carefully using random effects. Alternatively, if fixed effects vary at the level of the chick,  
82 then non-independence among clutches or mothers can be accounted for. Random  
83 effects typically represent some grouping variable (Breslow and Clayton 1993) and  
84 allow the estimation of variance in the response variable within and among these  
85 groups. This reduces the probability of false positives (Type I error rates) and false  
86 negatives (Type II error rates, e.g. Crawley 2013). Second, inferring the magnitude of  
87 variation within and among statistical clusters or hierarchical levels can be highly  
88 informative in its own right. In our bird example, understanding whether there is more  
89 variation in a focal trait among females within a population, rather than among  
90 populations, might be a central goal of the study.

91 LMMs are powerful yet complex tools. Software advances have made these tools  
92 accessible to the non-expert and have become relatively straightforward to fit in widely  
93 available statistical packages such as R (R Core Team 2016). Here we focus on the  
94 implementation of LMMs in R, although the majority of the techniques covered here can  
95 also be implemented in alternative packages including SAS (SAS Institute, Cary, NC) &  
96 SPSS (SPSS Inc., Chicago, IL). It should be noted however that due to different  
97 computational methods employed by different packages there may be differences in the  
98 model outputs generated. These differences will generally be subtle and the overall  
99 inferences drawn from the model outputs should be the same.

100 Despite this ease of implementation, the correct use of LMMs in the biological  
101 sciences is challenging for several reasons: i) they make additional assumptions about  
102 the data to those made in more standard statistical techniques such as general linear  
103 models (GLMs), and these assumptions are often violated (Bolker et al 2009); ii)  
104 interpreting model output correctly can be challenging, especially for the variance  
105 components of random effects (Bolker et al 2009; Zuur et al 2009); iii) model selection  
106 for LMMs presents a unique challenge, irrespective of model selection philosophy,  
107 because of biases in the performance of some tests (e.g. Wald tests, AIC comparisons)  
108 introduced by the presence of random effects (Vaida & Blanchard 2005; Dominicus et al  
109 2006; Bolker et al 2009). Collectively, these issues mean that the application of LMM  
110 techniques to biological problems can be risky and difficult for those that are unfamiliar  
4

111 with them. There have been several excellent papers in recent years on the use of  
112 generalized linear mixed effects models (GLMMs) in biology (Bolker et al 2009), the use  
113 of information theory and multi-model inference for studies involving LMMs (Grueber et  
114 al 2011), best practice for data exploration (Zuur et al 2009) and for conducting  
115 statistical analyses for complex datasets (Zuur & Ieno 2016; Kass et al 2016). At the  
116 interface of these excellent guides lies the theme of this paper: an updated guide for the  
117 uninitiated through the model fitting and model selection processes when using LMMs.  
118 A secondary but no less important aim of the paper is to bring together several key  
119 references on the topic of LMMs, and in doing so act as a portal into the primary  
120 literature that derives, describes and explains the complex modelling elements in more  
121 detail.

122 We provide a best practice guide covering the full analysis pipeline, from  
123 formulating hypotheses, specifying model structure and interpreting the resulting  
124 parameter estimates. The reader can digest the entire paper, or snack on each  
125 standalone section when required. First, we discuss the advantages and disadvantages  
126 of including both fixed and random effects in models. We then address issues of model  
127 specification, and choice of error structure and/or data transformation, a topic that has  
128 seen some debate in the literature (e.g. O'Hara & Kotze 2010; Ives 2015). We also  
129 address methods of model selection, and discuss the relative merits and potential  
130 pitfalls of using information theory (IT), AIC and multi-model inference in ecology and  
131 evolution. At all stages, we provide recommendations for the most sensible manner to  
132 proceed in different scenarios.

## 133 Understanding Fixed and Random Effects

134  
135 A key decision of the modelling process is specifying model predictors as fixed or  
136 random effects. Unfortunately, the distinction between the two is not always obvious,  
137 and is not helped by the presence of multiple, often confusing definitions in the literature  
138 (see Gelman and Hill 2007 p. 245). Absolute rules for how to classify something as a  
139 fixed or random effect generally are not useful because that decision can change

140 depending on the goals of the analysis (Gelman and Hill 2007). We can illustrate the  
141 difference between fitting something as a fixed (M1) or a random effect (M2) using a  
142 simple example of a researcher who takes measurements of mass from 100 animals  
143 from each of 5 different groups ( $n= 500$ ) with a goal of understanding differences among  
144 groups in mean mass. We use notation equivalent to fitting the proposed models in the  
145 statistical software *R* (R Core Team 2016), with the LMMs fitted using the R package  
146 *lme4* (Bates et al. 2015):

```
147  
148           M1 <- glm (mass ~ group)  
149           M2 <- lmer(mass ~ 1 + (1|group))
```

150  
151 Fitting 'group' as a fixed effect in model M1 assumes the 5 'group' means are all  
152 independent of one another, and share a common residual variance. Conversely, fitting  
153 group as a random intercept model in model M2 assumes that the 5 measured group  
154 means are only a subset of the realised possibilities drawn from a 'global' set of  
155 population means that follow a Normal distribution with its own mean ( $\mu_{\text{group}}$ , Fig. 1A)  
156 and variance ( $\sigma^2_{\text{group}}$ ). Therefore, LMMs model the variance hierarchically, estimating  
157 the processes that generate among-group variation in means, as well as variation within  
158 groups. Treating groups from a field survey as only a subset of the *possible* groups that  
159 could be sampled is quite intuitive, because there are likely many more groups (e.g.  
160 populations) of the study species in nature than the 5 the researcher measured.  
161 Conversely if one has designed an experiment to test the effect of three different  
162 temperature regimes on growth rate of plants, specifying temperature treatment as a  
163 fixed effect appears sensible because the experimenter has deliberately set the variable  
164 at a given value of interest. That is, there are no unmeasured groups with respect to  
165 that particular experimental design.

166 Estimating group means from a common distribution with known (estimated)  
167 variance has some useful properties, which we discuss below, and elaborate on the  
168 difference between fixed and random effects by using examples of the different ways  
169 random effects are used in the literature.

170

171 *Controlling for non-independence among data points*

172 This is one of the most common uses of a random effect. Complex biological data sets  
173 often contain nested and/or hierarchical structures such as repeat measurements from  
174 individuals within and across units of time. Random effects allow for the control of non-  
175 independence by constraining non-independent 'units' to have the same intercept  
176 and/or slope (Zuur et al 2009; Zuur & Ieno 2016). Fitting only random intercepts, or both  
177 random intercepts and slopes, will be decided by the goals of the analysis and the  
178 dependency structure of the data (Zuur & Ieno 2016). Fitting *only* a random intercept  
179 allows group means to vary, but assumes all groups have a common slope for a fitted  
180 covariate (fixed effect). Fitting random intercepts *and* slopes allows the slope of a  
181 predictor to vary based on a separate grouping variable. For example, one hypothesis  
182 might be that the probability of successful breeding for an animal is a function of its  
183 body mass. If we had measured animals from multiple sampling sites, we might wish to  
184 fit 'sampling site' as a random intercept, and estimate a common slope (change in  
185 breeding success) for body mass across all sampling sites by fitting it as a fixed effect:

186

```
187 M3 <- glmer(successful.breed ~ body.mass + (1|sample.site)
```

188

189 Conversely, we might wish to test the hypothesis that the strength of the effect (slope)  
190 of body mass on breeding success varies depending on the sampling location i.e. the  
191 change in breeding success for a 1 unit change in body mass is not consistent across  
192 groups (Figure 1B). Here, 'body mass' is specified as a random slope by moving it into  
193 the random effects structure:

194

```
195 M4 <- glmer(successful.breed ~ body.mass +  
196 (body.mass|sample.site)
```

197

198 Schielzeth & Forstmeier (2009) warn that constraining groups to share a common slope  
199 can inflate Type I and Type II errors. Consequently, Grueber et al (2011) recommend  
200 always fitting both random slopes and intercepts where possible. Whether this is  
201 feasible or not will depend on the data structure (see 'Costs to Fitting Random Effects'

7



202 section below). Figure 1 describes the differences between random intercept models  
203 and those also containing random slopes.

204 *Further reading: Zuur & Ieno (2016) shows examples of the difficulties in*  
205 *identifying the dependency structure of data and how to use flow charts / graphics to*  
206 *help decide model structure. Kery (2010, Ch 12) has an excellent demonstration of how*  
207 *to fit random slopes, and how model assumptions change depending on specification of*  
208 *a correlation between random slopes and intercepts or not. Schielzeth & Forstmeier*  
209 *(2009) and van de Pol & Wright (2009) are useful references for understanding the*  
210 *utility of random slope models.*

211

212 *Improving the accuracy of parameter estimation*

213 Random effect models use data from all the groups to estimate the mean and variance  
214 of the global distribution of group means. Assuming all group means are drawn from a  
215 common distribution causes the estimates of their means to drift towards the global  
216 mean  $\mu_{\text{group}}$ . This phenomenon, known as *shrinkage* (Gelman & Hill 2007; Kery 2010),  
217 can also lead to smaller and more precise standard errors around means. Shrinkage is  
218 strongest for groups with small sample sizes, as the paucity of within-group information  
219 to estimate the mean is counteracted by the model using data from other groups to  
220 improve the precision of the estimate. This 'partial pooling' of the estimates is a principal  
221 benefit of fitting something as a random effect (Gelman & Hill 2007). However, it can  
222 feel strange that group means should be shrunk towards the global mean, especially for  
223 researchers more used to treating sample means as independent fixed effects.

224 Accordingly, one issue is that variance estimates can be hugely imprecise when there  
225 are fewer than 5 levels of the random grouping variable (intercept or slope; see Harrison  
226 2015). However, thanks to the Central Limit Theorem, the assumption of Gaussian  
227 distribution of group means is usually a good one, and the benefits of hierarchical  
228 analysis will outweigh the apparent costs of shrinkage.

229

230 *Estimating variance components*

231 In some cases, the variation among groups will be of interest to ecologists. For  
232 example, imagine we had measured the clutch masses of 30 individual birds, each of  
8

233 which had produced 5 clutches (n=150). We might be interested in asking whether  
234 different females tend to produce consistently different clutch masses (high among-  
235 female variance for clutch mass). To do so, we might fit an intercept-only model with  
236 Clutch Mass as the response variable and a Gaussian error structure:

237

```
238     Model <- lmer(ClutchMass ~ 1 + (1|FemaleID)
```

239

240 By fitting individual 'FemaleID' as a random intercept term in the LMM, we estimate the  
241 among-female variance in our trait of interest. This model will also estimate the residual  
242 variance term, which we can use in conjunction with the among-female variance term to  
243 calculate an 'intra-class correlation coefficient' that measures individual repeatability in  
244 our trait (see Nakagawa & Schielzeth 2010). While differences among individuals can  
245 be obtained by fitting individual ID as a fixed effect, this uses a degree of freedom for  
246 each individual ID after the first, severely limiting model power, and does not benefit  
247 from increased estimation accuracy through shrinkage. More importantly, repeatability  
248 scores derived from variance components analysis can be compared across studies for  
249 the same trait, and even across traits in the same study. Variance component analysis  
250 is a powerful tool for partitioning variation in a focal trait among biologically interesting  
251 groups, and several more complex examples exist (see Nakagawa & Schielzeth 2010;  
252 Wilson et al 2010; Houslay & Wilson 2017). In particular, quantitative genetic studies  
253 rely on variance component analysis for estimating the heritability of traits such as body  
254 mass or size of secondary sexual characteristics (Wilson et al 2010). We recommend  
255 the tutorials in Wilson et al (2010) and Houslay & Wilson (2017) for a deeper  
256 understanding of the power and flexibility of variance component analysis.

257

258 *Making predictions for unmeasured groups*

259 Fixed effect estimates prevent us from making predictions for new groups because the  
260 model estimates are only relevant to groups in our dataset (e.g. Zuur et al 2009 p. 327).  
261 Conversely, we can use the estimate of the global distribution of population means to  
262 predict for the average group using the mean of the distribution  $\mu_{\text{group}}$  for a random  
263 effects model (see Fig. 1). We could also sample hypothetical groups from our random  
9

264 effect distribution, as we know its mean and SD (Zuur & Ieno 2016). Therefore, whether  
265 something is fitted as a fixed or random effect can depend on the goal of the analysis:  
266 are we only interested in the mean values for each group in our dataset, or do we wish  
267 to use our results to extend our predictions to new groups? Even if we do not want to  
268 predict to new groups, we might wish to fit something as a random effect to take  
269 advantage of the shrinkage effect and improved parameter estimation accuracy.

270

### 271 **Considerations When Fitting Random Effects**

272 Random effect models have several desirable properties (see above), but their use  
273 comes with some caveats. First, they are quite ‘data hungry’; requiring at least 5 ‘levels’  
274 (groups) for a random intercept term to achieve robust estimates of variance (Gelman &  
275 Hill 2007; Harrison 2015). With <5 levels, the mixed model may not be able to estimate  
276 the among-population variance accurately. In this case, the variance estimate will either  
277 collapse to zero, making the model equivalent to an ordinary GLM (Gelman & Hill 2007  
278 p. 275) or be non-zero but incorrect if the small number of groups that were sampled  
279 are not representative of true distribution of means (Harrison 2015). Second, models  
280 can be unstable if sample sizes across groups are highly unbalanced i.e. if some groups  
281 contain very few data. These issues are especially relevant to random slope models  
282 (Grueber et al 2011). Third, an important issue is the difficulty in deciding the  
283 “significance” or “importance” of variance among groups. The variance of a random  
284 effect is inevitably at least zero, but how big does it need to be to be considered of  
285 interest? Fitting a factor as a fixed effect provides a statement of the significance of  
286 differences (variation) among groups relatively easily. Testing differences among levels  
287 of a random effect is made much more difficult for frequentist analyses, though not so in  
288 a Bayesian framework (Kery 2010, see ‘*Testing Significance of Random Effects*’  
289 section). Finally, an issue that is not often addressed is that of mis-specification of  
290 random effects. GLMMs are powerful tools, but incorrectly parameterising the random  
291 effects in the model could yield model estimates that are as unreliable as ignoring the  
292 need for random effects altogether. An example would be failure to recognise non-  
293 independence caused by nested structures in the data e.g. multiple clutch measures  
294 from a single bird. A second example would be testing the significance of fixed effects at

10

295 the wrong 'level' of hierarchical models that ultimately leads to pseudoreplication and  
296 inflated Type I error rates. That is, if we take 10 measurements from each of 10 leaves  
297 to measure plant hormone concentration, even if we control for measurement non-  
298 independence with a random intercept for 'leaf ID', do we calculate our residual degrees  
299 of freedom at the data level (max n=100), or the grouping level (max n=10)?

300 *Further reading: Harrison (2015) shows how poor replication of the random*  
301 *intercept groups can give unstable model estimates. Zuur & Ieno (2016) discuss the*  
302 *importance of identifying dependency structures in the data.*

## 303 Deciding Model Structure for GLMMs

### 304 Choosing Error Structures and Link Functions

305 Linear models make various statistical assumptions, including additivity of the linear  
306 predictors, independence of errors, equal variance of errors (homoscedasticity) and  
307 Normality of errors (Gelman & Hill 2007 p. 46; Zuur et al 2009 p. 19). Ecologists often  
308 deal with response variables that violate these assumptions, and face several decisions  
309 about model specification to ensure models of such data are robust. The price for  
310 ignoring violation of these assumptions tends to be an inflated Type I error rate (Zuur et  
311 al 2010; Ives 2015). In some cases, however, transformation of the response variable  
312 may be required to ensure these assumptions are met. For example, an analytical goal  
313 may be to quantify differences in mean mass between males and females, but if the  
314 variance in mass for one sex is greater than the other, the assumption of homogeneity  
315 of variance is violated. Transformation of the data can remedy this (Zuur et al 2009),  
316 'mean-variance stabilising transformations' ensure the variance around the fitted mean  
317 of each group is similar, making the models more robust. Alternatively, modern  
318 statistical tools such as the 'varIdent' function in the R package *nlme* can allow one to  
319 explicitly model differences in variance between groups to avoid the need for data  
320 transformation.

321 *Further reading: Zuur et al (2010) provide a comprehensive guide on using data*  
322 *exploration techniques to check model assumptions, and give advice on*  
323 *transformations.*

324

325 For non-Gaussian data, our modelling choices become more complex. Non-  
326 Gaussian data structures include Poisson-distributed counts (number of eggs laid,  
327 number of parasites); binomial-distributed constrained counts (number of eggs that  
328 hatched in a clutch; prevalence of parasitic infection in a group of hosts) and Bernoulli-  
329 distributed binary traits (e.g. infected with a parasite or not). Gaussian models of these  
330 data would violate the assumptions of normality of errors and homogenous variance. To  
331 model these data, we have two initial choices: i) we can apply a transformation to our  
332 non-Gaussian response to 'make it' approximately Gaussian, and then use a Gaussian  
333 model; or ii) we can apply a GL(M)M and specify the appropriate error distribution and  
334 link function. The link function takes into account the (assumed) empirical distribution of  
335 our data by transformation of the linear predictor within the model. It is critical to note  
336 that transformation of the raw response variable is not equivalent to using a link function  
337 to apply a transformation in the model. Data-transformation applies the transformation  
338 to the raw response, whilst using a link function transforms the fitted mean (the linear  
339 predictor). That is, *the mean of a log-transformed response (using a data*  
340 *transformation) is not identical to the logarithm of a fitted mean (using a link function).*

341 The issue of transforming non-Gaussian data to fit Gaussian models to them is  
342 contentious. For example, arcsin square-root transformation of proportion data was  
343 once extremely common, but recent work has shown it to be unreliable at detecting real  
344 effects (Warton & Hui 2011). Both logit-transformation (for proportional data) and  
345 Binomial GLMMs (for binary response variables) have been shown to be more robust  
346 (Warton & Hui 2011). O'Hara & Kotze (2010) argued that log-transformation of count  
347 data performed well in only a small number of circumstances (low dispersion, high  
348 mean counts), which are unlikely to be applicable to ecological datasets. However, Ives  
349 (2015) recently countered these assumptions with evidence that transformed count data  
350 analysed using LMMs can often outperform Poisson GLMMs. We do not make a case  
351 for either here, but acknowledge the fact that there is unlikely to be a universally best  
12

352 approach; each method will have its own strengths and weakness depending on the  
353 properties of the data (O'Hara & Kotze 2010). Checking the assumptions of the LMM or  
354 GLMM is an essential step. An issue with transformations of non-Gaussian data is  
355 having to deal with zeroes as special cases (e.g. you can't log transform a 0), so  
356 researchers often add a small constant to all data to make the transformation work, a  
357 practice that has been criticised (O'Hara & Kotze 2010). GLMMs remove the need for  
358 these 'adjustments' of the data. The important point here is that transformations change  
359 the entire relationship between Y and X (Zuur et al 2009), but different transformations  
360 do this to different extents and it may be impossible to know which transformation is  
361 best without performing simulations to test the efficacy of each (Warton & Hui 2011;  
362 Ives 2015).

363 *Further reading: Crawley (2013 Ch 13) gives a broad introduction to the various error*  
364 *structures and link functions available in the R statistical framework. O'Hara & Kotze*  
365 *(2010); Ives (2015) and Warton et al (2016) argue the relative merits of GLMs vs log-*  
366 *transformation of count data; Warton & Hui (2011) address the utility of logit-*  
367 *transformation of proportion data compared to arcsin square-root transformation.*

368

### 369 **Choosing Random Effects I: Crossed or Nested?**

370 A common issue that causes confusion is this issue of specifying random effects as  
371 either 'crossed' or 'nested'. In reality, the way you specify your random effects will be  
372 determined by your experimental or sampling design (Schielzeth & Nakagawa 2013). A  
373 simple example can illustrate the difference. Imagine a researcher was interested in  
374 understanding the factors affecting the clutch mass of a passerine bird. They have a  
375 study population spread across 5 separate woodlands, each containing 30 nest boxes.  
376 Every week during breeding they measure the foraging rate of females at feeders, and  
377 measure their subsequent clutch mass. Some females have multiple clutches in a  
378 season and contribute multiple data points. Here, female ID is said to be *nested within*  
379 *woodland*: each woodland contains multiple females unique to that woodland (that  
380 never move among woodlands). The nested random effect controls for the fact that i)  
381 clutches from the same female are not independent, and ii) females from the same  
13

382 woodland may have clutch masses more similar to one another than to females from  
383 other woodlands

384  
385 `Clutch Mass ~ Foraging Rate + (1|Woodland/Female ID)`

386  
387 Now imagine that this is a long-term study, and the researcher returns every year for 5  
388 years to continue with measurements. Here it is appropriate fit year as a *crossed*  
389 random effect, because every woodland appears multiple times in every year of the  
390 dataset, and females that survive from one year to the next will also appear in multiple  
391 years.

392  
393 `Clutch Mass ~ Foraging Rate + (1|Woodland/Female ID)+ (1|Year)`

394  
395 Understanding whether your experimental/sampling design calls for nested or crossed  
396 random effects is not always straightforward, but it can help to visualise experimental  
397 design by drawing it (see Schielzeth and Nakagawa 2013 Fig. 1), or tabulating your  
398 observations by these grouping factors (e.g. with the *'table'* command in R) to identify  
399 how your data are distributed. Finally, we caution that whether two factors are nested or  
400 crossed affects the ability of GLMMs to estimate the interaction variance between those  
401 two groups on the outcome variable. Crossed factors can accurately estimate the  
402 interaction variance between the two, whereas nested factors automatically pool the  
403 interaction variance in the second (nested) factor (Schielzeth and Nakagawa 2013). We  
404 do not expand on this important issue here but direct the reader to Schielzeth and  
405 Nakagawa 2013 for an excellent treatment of the topic.

## 406 **Choosing Random Effects II: Random Slopes for Continuous Variables**

407 Fitting random slope models in ecology is not very common. Often, researchers fit  
408 random intercepts to control for non-independence among measurements of a statistical  
409 group (e.g. birds within a woodland), but allow a continuous variable to have a common  
410 slope across all experimental units. Schielzeth & Forstmeier (2009) argue that including  
411 random slopes controls Type I error rate for continuous predictors (yields more accurate  
14

412 p values), but also give more power to detect among individual variation. Barr et al  
413 (2013) argue that researchers should fit the maximal random effects structure possible  
414 for the data. That is, if there are four continuous predictors under consideration, all four  
415 should be allowed to have random slopes. However, we believe this is unrealistic  
416 because random slope models require large numbers of data to estimate variances and  
417 covariances accurately (Bates et al 2015). Ecological datasets can often struggle to  
418 estimate a single random slope, diagnosed by a perfect correlation (1 or -1) between  
419 random intercepts and slopes (Bates et al 2015). Therefore, the approach of fitting the  
420 'maximal' complexity of random effects structure (Barr et al 2013) is perhaps better  
421 phrased as fitting the most complex mixed effects structure allowed by your data (Bates  
422 et al 2015), which may mean no random slopes at all. If fitting a random slope model,  
423 always inspect the correlation coefficient between the intercepts and slopes in the  
424 variance/covariance summary returned by packages like *lme4* to look for evidence of  
425 perfect correlations, indicative of insufficient data to estimate the model.

426 *Further Reading: Forstmeier and Schielzeth (2009) is essential reading for*  
427 *understanding how random slopes control Type I error rate, and Bates et al (2015)*  
428 *gives sound advice on how to iteratively determine optimal complexity of random effect*  
429 *structure.*

### 430 **Choosing Fixed Effect Predictors and Interactions**

431 One of the most important decisions during the modelling process is deciding which  
432 predictors and interactions to include in models. Best practice demands that each model  
433 should represent a specific *a priori* hypothesis concerning the drivers of patterns in data  
434 (Burnham & Anderson 2002; Forstmeier & Schielzeth 2011), allowing the assessment of  
435 the relative support for these hypotheses in the data irrespective of model selection  
436 philosophy. The definition of "hypothesis" must be broadened from the strict pairing of  
437 null and alternative that is classically drilled into young pupils of statistics and  
438 experimental design. Frequentist approaches to statistical modelling still work with  
439 nested pairs of hypotheses. Information theorists work with whole sets of competing  
440 hypotheses. Bayesian modellers are comfortable with the idea that every possible  
441 parameter estimate is a hypothesis in its own right. But these epistemological

15



442 differences do not really help to solve the problem of “which” predictors should be  
443 considered valid members of the full set to be used in a statistical modelling exercise. It  
444 is therefore often unclear how best to design the most complex model, often referred to  
445 as the *maximal model* (which contains all factors, interactions and covariates that might  
446 be of any interest, Crawley 2013) or as the *global model* (a highly parameterized model  
447 containing the variables and associated parameters thought to be important of the  
448 problem at hand, Burnham & Anderson 2002; Grueber et al 2011). We shall use the  
449 latter term here for consistency with terminology used in information-theory (Grueber et  
450 al 2011).

451         Deciding which terms to include in the model requires careful and rigorous *a*  
452 *priori* consideration of the system under study. This may appear obvious; however  
453 diverse authors have noticed a lack of careful thinking when selecting variables for  
454 inclusion in a model (Peters 1991, Chatfield 1995, Burnham & Anderson 2002). Lack of  
455 *a priori* consideration, of what models represent, distinguishes rigorous hypothesis  
456 testing from ‘fishing expeditions’ that seek significant predictors among a large group of  
457 contenders. Ideally, the global model should be carefully constructed using the  
458 researchers’ knowledge and understanding of the system such that only predictors likely  
459 to be pertinent to the problem at hand are included, rather than including all the data the  
460 researcher has collected and/or has available. This is a pertinent issue in the age of ‘big  
461 data’, where researchers are often overwhelmed with predictors and risk skipping the  
462 important step of *a priori* hypothesis design. In practice, for peer reviewers it is easy to  
463 distinguish fishing expeditions from *a priori* hypothesis sets based on the evidence base  
464 presented in introductory sections of research outputs.

465

## 466 **How Complex Should My Global Model Be?**

467         The complexity of the global model will likely be a trade-off between the number  
468 of measured observations (the  $n$  of the study) and the proposed hypotheses about how  
469 the measured variables affect the outcome (response) variable. Lack of careful  
470 consideration of the parameters to be estimated can result in overparameterised  
471 models, where there are insufficient data to estimate coefficients robustly (Southwood &  
16

472 Henderson 2000, Quinn & Keough 2002, Crawley 2013). In simple GLMs,  
473 overparameterisation results in a rapid decline in (or absence of) degrees of freedom  
474 with which to estimate residual error. Detection of overparameterisation in LMMs can be  
475 more difficult because each random effect uses only a single degree of freedom,  
476 however the estimation of variance among small numbers of groups can be numerically  
477 unstable. Unfortunately, it is common practice to fit a global model that is simply as  
478 complex as possible, irrespective of what that model actually represents; that is a  
479 dataset containing  $k$  predictors yields a model containing a  $k$ -way interaction among all  
480 predictors and simplify from there (Crawley 2013). This approach is flawed for two  
481 reasons. First, this practice encourages fitting biologically-unfeasible models containing  
482 nonsensical interactions. It should be possible to draw and/or visualise what the fitted  
483 model 'looks like' for various combinations of predictors – being unable to draw the  
484 expected fitted lines of a 3-way interaction means refraining from fitting a model  
485 containing one. Second, using this approach makes it very easy to fit a model too  
486 complex for the data. At best, the model will fail to converge, thus preventing inference.  
487 At worst, the model will “work”, risking false inference. Guidelines for the ideal ratio of  
488 data points ( $n$ ) to estimated parameters ( $k$ ) vary widely (see Forstmeier & Schielzeth  
489 2011). Crawley (2013) suggests a minimum  $n/k$  of 3, though we argue this is very low  
490 and that an  $n/k$  of 10 is more conservative. A 'simple' model containing a 3-way  
491 interaction between continuous predictors and a single random intercept needs to  
492 estimate 8 parameters, so requires a dataset of a *minimum*  $n$  of 80. Interactions can be  
493 especially demanding, as fitting interactions between a multi-level factor and a  
494 continuous predictor can result in poor sample sizes for specific treatment combinations  
495 even if the total  $n$  is quite large (Zuur et al 2010), which will lead to unreliable model  
496 estimates.

497 *Grueber et al (2011) show an excellent worked example of a case where the*  
498 *most complex model is biologically feasible and well-reasoned, containing only one 2-*  
499 *way interaction. Nakagawa and Foster (2004) discuss the use of power analyses, which*  
500 *will be useful in determining the appropriate  $n/k$  ratio for a given system.*

501

502 *Assessing Predictor Collinearity*

17

503 With the desired set of predictors identified, it is wise to check for collinearity among  
504 predictor variables. Collinearity among predictors can cause several problems in model  
505 interpretation because those predictors explain some of the same variance in the  
506 response variable, and their effects cannot be estimated independently (Quinn and  
507 Keough. 2002; Graham 2003): First, it can cause model convergence issues as models  
508 struggle to partition variance between predictor variables. Second, positively correlated  
509 variables can have negatively correlated regression coefficients, as the marginal effect  
510 of one is estimated, given the effect of the other, leading to incorrect interpretations of  
511 the direction of effects (Figure 2). Third, collinearity can inflate standard errors of  
512 coefficient estimates and make 'true' effects harder to detect (Zuur et al 2010). Finally,  
513 collinearity can affect the accuracy of model averaged parameter estimates during  
514 multi-model inference (Freckleton 2011; Cade 2015). Examples of collinear variables  
515 include climatic data such as temperature and rainfall, and morphometric data such as  
516 body length and mass. Collinearity can be detected in several ways, including creating  
517 correlation matrices between raw explanatory variables, with values  $>0.7$  suggesting  
518 both should not be used in the same model (Dormann et al. 2013); or calculating the  
519 variance inflation factor (VIF) of each predictor that is a candidate for inclusion in a  
520 model (details in Zuur et al 2010) and dropping variables with a VIF higher than a  
521 certain value (e.g. 3; Zuur et al 2010, or 10, Quinn & Keogh 2002). One problem with  
522 these methods though is that they rely on a user-selected choice of threshold of either  
523 the correlation coefficient or the VIF, and use of more stringent (lower) is probably  
524 sensible. Some argue that one should always prefer inspection of VIF values over  
525 correlation coefficients of raw predictors because strong multicollinearity can be hard to  
526 detect with the latter. When collinearity is detected, researchers can either select one  
527 variable as representative of multiple collinear variables (Austin 2002), ideally using  
528 biological knowledge/ reasoning to select the most meaningful variable (Zuur et al  
529 2010); or conduct a dimension-reduction analysis (e.g. Principal Components Analysis;  
530 James & McCullough 1990), leaving a single variable that accounts for most of the  
531 shared variance among the correlated variables. Both approaches will only be  
532 applicable if it is possible to group explanatory variables by common features, thereby  
533 effectively creating broader, but still meaningful explanatory categories. For instance, by  
18

534 using mass and body length metrics to create a 'scaled mass index' representative of  
535 body size (Peig & Green 2009).

536

### 537 *Standardising and Centering Predictors*

538 Transformations of predictor variables are common, and can improve model  
539 performance and interpretability (Gelman & Hill 2007). Two common transformations for  
540 continuous predictors are i) predictor centering, the mean of predictor  $x$  is subtracted  
541 from every value in  $x$ , giving a variable with mean 0 and SD on the original scale of  $x$ ;  
542 and ii) predictor standardising, where  $x$  is centred and then divided by the SD of  $x$ ,  
543 giving a variable with mean 0 and SD 1. Rescaling the mean of predictors containing  
544 large values (e.g. rainfall measured in thousands of mm) through  
545 centering/standardising will often solve convergence problems, in part because the  
546 estimation of intercepts is brought into the main body of the data themselves. Both  
547 approaches also remove the correlation between main effects and their interactions,  
548 making main effects interpretable when models also contain interactions (Schielzeth  
549 2010). Note that this collinearity among coefficients is distinct from collinearity between  
550 two separate predictors (see above). Centering and standardising by the mean of a  
551 variable changes the interpretation of the model intercept to the value of the outcome  
552 expected when  $x$  is at its mean value. Standardising further adjusts the interpretation of  
553 the coefficient (slope) for  $x$  in the model to the change in the outcome variable for a 1  
554 SD change in the value of  $x$ . Scaling is therefore a useful, indeed recommended, tool to  
555 improve the stability of models and likelihood of model convergence, and the accuracy  
556 of parameter estimates, but care must be taken in the interpretation and graphical  
557 representation of outcomes.

558 *Further reading: Schielzeth (2010) provides an excellent reference to the*  
559 *advantages of centering and standardising predictors. Gelman (2008) provides strong*  
560 *arguments for standardising continuous variables by 2 SDs when binary predictors are*  
561 *in the model. Gelman & Hill (2007 p. 56, 434) discuss the utility of centering by values*  
562 *other than the mean.*

563

## 564 Quantifying GLMM Fit and Performance

565 Once a global model is specified, it is vital to quantify model fit and report these metrics  
566 in the manuscript. The global model is considered the best candidate for assessing fit  
567 statistics such as overdispersion (Burnham & Anderson 2002). Information criteria  
568 scores should not be used as a proxy for model fit, because a large difference in AIC  
569 between the top and null models is not evidence of a good fit. AIC tells us nothing about  
570 whether the basic distributional and structural assumptions of the model have been  
571 violated. Similarly a high  $R^2$  value is in itself only a test of the magnitude of model fit and  
572 not an adequate surrogate for proper model checks. Just because a model has a high  
573  $R^2$  value does not mean it will pass checks for assumptions such as homogeneity of  
574 variance. We strongly encourage researchers to view *model fit* and *model adequacy* as  
575 two separate but equally important traits that must be assessed and reported. Model fit  
576 can be poor for several reasons, including the presence of overdispersion, failing to  
577 include interactions among predictors, failing to account for non-linear effects of the  
578 predictors on the response, or specifying a sub-optimal error structure and/or link  
579 function. Here we discuss some key metrics of fit and adequacy that should be  
580 considered.

581

### 582 *Inspection of Residuals and Linear Model Assumptions*

583 Best practice is to examine plots of fitted values vs residuals for the entire model, as  
584 well as model residuals versus all explanatory variables to look for patterns (Zuur et al  
585 2010; Zuur & Ieno 2016). In addition, there are further model checks specific to mixed  
586 models. First, inspect fitted values versus residuals for each grouping level of a random  
587 intercept factor (Zuur et al 2009). This will often prove dissatisfying if there are few  
588 data/residuals per group, however this in itself is a warning flag that the assumptions of  
589 the model might be based on weak foundation. Note that for the GLMMs it is wise to  
590 use normalised/Pearson residual when looking for patterns as they account for the  
591 mean-variance relationship of generalized models (Zuur et al 2009). Another feature of  
592 fit that is very rarely tested for in (G)LMMs is the assumption of normality of deviations  
593 of the conditional means of the random effects from the global intercept. Just as a

594 quantile-quantile (QQ) plot of linear model residuals should show points falling along a  
595 straight line (e.g. Crawley 2013), so should a QQ plot of the random effect means  
596 (Schielzeth & Nakagawa 2013).

597 *Further reading: Zuur et al (2010) given an excellent overview of the assumptions of*  
598 *linear models and how to test for their violation. See also Gelman & Hill (2007 p. 45).*

599 *The R package 'sjPlot' (Lüdecke 2017) has built in functions for several LMM*  
600 *diagnostics, including random effect QQ plots. Zuur et al (2009) provides a vast*  
601 *selection of model diagnostic techniques for a host of model types, including GLS,*  
602 *GLMMs and GAMMS.*

603

604 *Overdispersion*

605 Models with a Gaussian (Normal) error structure do not require adjustment for  
606 overdispersion, as Gaussian models do not assume a specific mean-variance  
607 relationship. For generalized mixed models (GLMMs) however (e.g. Poisson, Binomial),  
608 the variance of the data can be greater than predicted by the error structure of the  
609 model (e.g. Hilbe 2011). Overdispersion can be caused by several processes  
610 influencing data, including zero-inflation, aggregation (non-independence) among  
611 counts, or both (Zuur et al 2009). The presence of overdispersion in a model suggests it  
612 is a bad fit, and standard errors of estimates will likely be biased unless overdispersion  
613 is accounted for (e.g. Harrison 2014). The use of canonical binomial and Poisson error  
614 structures, when residuals are overdispersed, tends to result in Type I errors because  
615 standard errors are underestimated. Adding an observation-level random effect (OLRE)  
616 to overdispersed Poisson or Binomial models can model the overdispersion and give  
617 more accurate estimates standard errors (Harrison 2014; 2015). However, OLRE  
618 models may yield inferior fit and/or biased parameter estimates compared to models  
619 using compound probability distributions such as the Negative-Binomial for count data  
620 (Hilbe 2011; Harrison 2014) or Beta-Binomial for proportion data (Harrison 2015), and  
621 so it is good practice to assess the relative fit of both types of model using AIC before  
622 proceeding (e.g. Zuur et al 2009). Researchers very rarely report the overdispersion  
623 statistic (but see Elston et al 2001), but it should be made a matter of routine. See

624 'Assessing Model Fit Through Simulation' Section for advice on how to quantify and  
625 model overdispersion.

626 *Further reading: Crawley (2013 page 580-581) gives an elegant demonstration of*  
627 *how failing to account for overdispersion leads to artificially small standard errors and*  
628 *spurious significance of variables. Harrison (2014) quantifies the ability of OLRE to cope*  
629 *with overdispersion in Poisson models. Harrison (2015) compares Beta-Binomial and*  
630 *OLRE models for overdispersed proportion data.*

631

632  $R^2$

633 In a linear modelling context,  $R^2$  gives a measure of the proportion of explained variance  
634 in the model, and is an intuitive metric for assessing model fit. Unfortunately, the issue  
635 of calculating  $R^2$  for (G)LMMs is particularly contentious; whereas residual variance can  
636 easily be estimated for a simple linear model with no random effects and a Normal error  
637 structure, this is not the case for (G)LMMs. In fact, two issues exist with generalising  $R^2$   
638 measures to (G)LMMs: i) for generalised models containing non-Normal error  
639 structures, it is not clear how to calculate the residual variance term on which the  $R^2$   
640 term is dependent; and ii) for mixed effects models, which are hierarchical in nature and  
641 contain error (unexplained variance) at each of these levels, it is uncertain which level to  
642 use to calculate a residual error term (Nakagawa & Schielzeth 2013). Diverse methods  
643 have been proposed to account for this coefficient in GLMMs, including so-called  
644 'pseudo- $r^2$ ' measures of explained variance (e.g. Nagelkerke 1991, Cox & Snell 1989),  
645 but their performance is often unstable for mixed models and can return negative values  
646 (Nakagawa & Schielzeth 2013). Gelman & Pardoe (2006) derived a measure of  $R^2$  that  
647 accounts for the hierarchical nature of LMMs and gives a measure for both group and  
648 unit level regressions (see also Gelman & Hill 2007 p. 474), but it was developed for a  
649 Bayesian framework and a frequentist analogue does not appear to be widely  
650 implemented. The method that has gained the most support over recent years is that of  
651 Nakagawa & Schielzeth (2013).

652 The strength of the Nakagawa & Schielzeth (2013) method for GLMMs is that it  
653 returns two complimentary  $R^2$  values: the marginal  $R^2$  encompassing variance  
654 explained by only the fixed effects, and the conditional  $R^2$  comprising variance  
22

655 explained by both fixed and random effects i.e. the variance explained by the whole  
656 model (Nakagawa & Schielzeth 2013). Ideally, both should be reported in publications  
657 as they provide different information; which one is more 'useful' may depend on the  
658 rationale for specifying random effects in the first instance. Recently, Nakagawa,  
659 Johnson & Schielzeth (2017) expanded their  $R^2$  method to handle models with  
660 compound probability distributions like the Negative Binomial error family. Note that  
661 when observation-level random effects are included (see 'Overdispersion' section  
662 above), the conditional  $R^2$  becomes less useful as a measure of explained variance  
663 because it includes the extra-parametric dispersion being modelled, but has no  
664 predictive power (Harrison 2014).

665 *Further reading: Nakagawa & Schielzeth (2013) provide an excellent and*  
666 *accessible description of the problems with, and solutions to, generalising  $R^2$  metrics to*  
667 *GLMMs. The Nakagawa & Schielzeth (2013)  $R^2$  functions have been incorporated into*  
668 *several packages, including 'MuMIn' (Bartoń 2016) and 'piecewiseSEM' (Lefcheck*  
669 *2015), and Johnson (2014) has developed an extension of the functions for random*  
670 *slope models. See Harrison (2014) for a cautionary tale of how the GLMM  $R^2$  functions*  
671 *are artificially inflated for overdispersed models.*

672  
673

#### 674 *Stability of Variance Components and Testing Significance of Random Effects*

675 When models are too complex relative to the amount of data available, GLMM variance  
676 components can collapse to zero (they cannot be negative). This is not a problem *per*  
677 *se*, but it's important to acknowledge that in this case the model is equivalent to a  
678 standard GLM. Reducing model complexity by removing interactions will often allow  
679 random effects variance component estimates to become  $>0$ , but this is problematic if  
680 quantifying the interaction is the primary goal of the study. REML (restricted maximum  
681 likelihood) should be used for estimating variance components of random effects in  
682 Gaussian GLMMs as it produces less biased estimates compared to ML (maximum  
683 likelihood) (Bolker et al 2009). However, when comparing two models with the same  
684 random structure but different fixed effects, ML estimation cannot easily be avoided.  
685 The RLRsim package (Scheipl, 2016) can be used to calculate restricted likelihood ratio

23



686 tests for variance components in mixed and additive models. Crucially, when testing the  
687 significance of a variance component we are ‘testing on the boundary’ (Bolker et al  
688 2009). That is the null hypothesis for random effects ( $\sigma=0$ ) is at the boundary of its  
689 possible range (it has to be  $\geq 0$ ), meaning p-values from a likelihood ratio test are  
690 inaccurate. Dividing p values by 2 for tests of single variance components provides an  
691 approximation to remedy this problem (Verbenke & Molenberghs, 2000).

692 Finally, estimating degrees of freedom for tests of random effects using Wald, t  
693 or F tests or AICc is difficult, as a random effect can theoretically use anywhere  
694 between 1 and  $N - 1$  df (where N is the number of random-effect levels) (Bolker et al.  
695 2009). Adequate F and P values can be calculated using Satterthwaite (1946)  
696 approximations to determine denominator degrees of freedom implemented in the  
697 package ‘lmerTest’ (Kuznetzova et al. 2014, see further details in section ‘Model  
698 Selection and Multi-Model Inference’ below).

699

#### 700 *Assessing Model Fit through Simulation*

701 Simulation is a powerful tool for assessing model fit (Gelman & Hill 2007; Kery 2010;  
702 Zuur & Ieno 2016), but is rarely used. The premise here is simple: when simulating a  
703 dataset from a given set of parameter estimates (a model), the fit of the model to those  
704 *simulated* ‘ideal’ data should be comparable to the model’s fit to the real data (Kery  
705 2010). Each iteration yields a simulated dataset that allows calculation of a statistic of  
706 interest such as the sum of squared residuals (Kery 2010), the overdispersion statistic  
707 (Harrison 2014) or the percentage of zeroes for a Poisson model (Zuur & Ieno 2016). If  
708 the model is a good fit, after a sufficiently large number of iterations (e.g. 10,000) the  
709 distribution of this test statistic should encompass the observed statistic in the real data.  
710 Significant deviations outside of that distribution indicate the model is a poor fit (Kery  
711 2010). Figure 3 shows an example of using simulation to assess the fit of a Poisson  
712 GLMM. After fitting a GLMM to count data, we may wish to check for overdispersion  
713 and/or zero-inflation, the presence of which might suggest we need to adjust our  
714 modelling strategy. Simulating 10,000 datasets from our model reveals that the  
715 proportion of zeroes in our real data is comparable to simulated expectation (Figure 3A).  
716 Conversely, simulating 1000 datasets and refitting our model to each dataset, we see

717 that the sum of the squared Pearson residuals for the real data is far larger than  
718 simulated expectation (Figure 3B), giving evidence of overdispersion (Harrison 2014).  
719 We can use the simulated frequency distribution of this test statistic to derive a mean  
720 and 95% confidence interval for the overdispersion by calculating the ratio of our test  
721 statistic to the simulated values (Harrison 2014). The dispersion statistic for our model is  
722 3.16 [95% CI 2.77 – 3.59]. Thus, simulations have allowed us to conclude that our  
723 model is overdispersed, but that this overdispersion is not due to zero-inflation. All R  
724 code for reproducing these simulations is provided in Online Supplementary Material.

725 *Further reading: The R package ‘SQuID’ (Allegue et al 2017) provides a highly*  
726 *flexible simulation tool for learning about, and exploring the performance of, GLMMs.*  
727 *Rykiel (1996) discusses the need for validation of models in ecology.*

728

#### 729 *Dealing with missing data*

730 Often when collecting ecological data it is not always possible to measure all of the  
731 predictors of interest for every measurement of the dependant variable. Such missing  
732 data is a common feature of ecological datasets, however the impacts of this have  
733 seldom been considered in the literature (Nakagawa & Freckleton 2011). Incomplete  
734 data is often dealt with by deleting data point with missing predictor data (Nakagawa &  
735 Freckleton 2008), although this may results in biased parameter estimates and reduces  
736 statistical power (Nakagawa & Freckleton 2008). Nakagawa & Freckleton (2011)  
737 recommend multiple imputation (MI) as a mechanism for handling missing data, and  
738 highlight the ability of this technique for more accurate estimates, particularly for IT-AIC  
739 approaches.

740 *Further reading: See Nakagawa & Freckleton (2008) for a review on the risks of*  
741 *ignoring incomplete data. Nakagawa & Freckleton (2011) demonstrate the effects of*  
742 *missing data during model selection procedures, and provide an overview of R*  
743 *packages available for MI.*

## 744 Model Selection and Multi-Model Inference

745 Several methods of model selection are available once there is a robust global model  
746 that satisfies standard assumptions of error structure and hierarchical independence  
747 (Johnson & Omland 2004). We discuss the relative merits of each approach briefly  
748 here, before expanding on the use of information-theory and multi-model inference in  
749 ecology. We note that these discussions are not meant to be exhaustive comparisons,  
750 and we encourage the reader to delve into the references provided for a comprehensive  
751 picture of the arguments for and against each approach.

752

### 753 *Stepwise Selection, Likelihood Ratio Tests and P values*

754 A common approach to model selection is the comparison of a candidate model  
755 containing a term of interest to the corresponding 'null' model lacking that term, using a  
756 p value from a likelihood ratio test (LRT), referred to as null-hypothesis significance  
757 testing (NHST; Nickerson 2000). Stepwise deletion involves using the NHST framework  
758 to drop terms sequentially from the global model, and arrive at a 'minimal adequate  
759 model' (MAM) containing only significant predictors (see Crawley 2013). NHST and  
760 stepwise deletion have come under heavy criticism; they can overestimate the effect  
761 size of 'significant' predictors (Whittingham et al 2006; Forstmeier & Schielzeth 2011)  
762 and force the researcher to focus on a single best model as if it were the only  
763 combination of predictors with support in the data. Although we strive for simplicity and  
764 parsimony, this assumption is not reasonable in complex ecological systems (e.g.  
765 Burnham, Anderson & Huyvaert 2011). It is common to present the MAM as if it arose  
766 from a single *a priori* hypothesis, when in fact arriving at the MAM required multiple  
767 significance tests (Whittingham et al 2006; Forstmeier & Schielzeth 2011). This cryptic  
768 multiple testing can lead to hugely inflated Type I errors (Forstmeier & Schielzeth 2011).  
769 Perhaps most importantly, LRT can be unreliable for fixed effects in GLMMs unless both  
770 total sample size and replication of the random effect terms is high (see Bolker et al  
771 2009 and references therein), conditions which are often not satisfied for most  
772 ecological datasets. However, there are still cases where NHST may be the most  
773 appropriate tool for inference (Murtaugh 2014). For example, in controlled experimental

26

774 studies a researcher may wish to test the effect of a limited number of treatments and  
775 support estimates of effect sizes with statements of statistical significance using model  
776 simplification (Mundry 2011). Importantly, Murtaugh (2009) found that the predictive  
777 ability of models assessed using NHST was comparable to those selected using  
778 information-theoretic approaches (see below), suggesting that NHST remains a valid  
779 tool for inference despite strong criticism (see also Murtaugh 2014). Our advice is that  
780 NHST remains an important tool for analyses of experiments and for inferential surveys  
781 with small numbers of well-justified *a priori* hypotheses and with uncorrelated (or weakly  
782 correlated) predictors.

783 *Further reading: See Murtaugh's (2014) excellent 'in Defense of P values;', as*  
784 *well as the other papers on the topic in the same special issue of Ecology. Stephens et*  
785 *al (2005) & Mundry (2011) argue the case for NHST under certain circumstances such*  
786 *as well-designed experiments. Halsey et al (2015) discuss the wider issues of the*  
787 *reliability of p values relative to sample size.*

788

#### 789 *Information-Theory and Multi-Model Inference*

790 Unlike NHST, which leads to a focus on a single best model, model selection using  
791 information theoretic (IT) approaches allows the degree of support in the data for  
792 several competing models to be ranked using metrics such as Akaike's Information  
793 Criterion (AIC). Information criteria attempt to quantify the Kullback-Leibler distance  
794 (KLD), a measure of the relative amount of information lost when a given model  
795 approximates the true data-generating process. Thus, relative difference among models  
796 in AIC should be representative in relative differences in KLD, and the model with the  
797 lowest AIC should lose the least information and be the best model in that it optimises  
798 the trade-off between fit and complexity (e.g. Richards 2008). A key strength of the IT  
799 approach is that it accounts for 'model selection uncertainty', the idea that several  
800 competing models may all fit the data similarly (Burnham & Anderson 2002; Burnham,  
801 Anderson & Huyvaert 2011). This is particularly useful when competing models share  
802 equal "complexity" (i.e. number of predictors, or number of residual degrees of  
803 freedom): in such situations, NHST is impossible because there is no "null". Where  
804 several models have similar support in the data, inference can be made from all models

27

805 using model-averaging (Burnham & Anderson 2002; Johnson & Omand 2004; Grueber  
806 et al 2011). Model averaging incorporates uncertainty by weighting the parameter  
807 estimate of a model by that model's Akaike weight (often referred to as the probability of  
808 that model being the best Kullback-Leibler model given the data, but see Richards  
809 2005). Multi-model inference places a strong emphasis on *a priori* formulation of  
810 hypotheses (Burnham & Anderson 2002; Dochterman & Jenkins 2011; Lindberg et al  
811 2015), and model-averaged parameter estimates arising from multi-model inference are  
812 thought to lead to more robust conclusions about the biological systems compared to  
813 NHST (Johnson & Omland 2004, but see Richards et al 2011). These strengths over  
814 NHST have meant that the use of IT approaches in ecology and evolution has grown  
815 rapidly in recent years (Lindberg et al 2015; Barker & Link 2015; Cade 2015). We do not  
816 expand on the specific details of the difference between NHST and IT here, but point  
817 the reader to some excellent reference on the topic. Instead, we use this section to  
818 highlight recent empirical developments in the best practice methods for the application  
819 of IT in ecology and evolution.

820 *Further reading: Grueber et al (2011) and Symonds & Moussalli (2011) give a*  
821 *broad overview of multi-model inference in ecology, and provide a worked model*  
822 *selection exercise. Heygi & Garamszegi (2011) provide a detailed comparison of IT and*  
823 *NHST approaches. Burnham, Anderson & Huyvaert (2011) demonstrate how AIC*  
824 *approximates Kullback-Leibler information and provide some excellent guides for the*  
825 *best practice of applying IT methods to biological datasets. Vaida & Blanchard (2005)*  
826 *provide details on AIC should be implemented for the analysis of clustered data.*

827

### 828 *Global Model Reporting*

829 Because stepwise deletion can cause biased effect sizes, presenting means and SEs of  
830 parameters from the global model should be more robust, especially when the n/k ratio  
831 is low (Forstmeier & Schielzeth 2011). An alternative approach to NHST is to perform  
832 'full model tests' (comparing the global model to an intercept only model) before  
833 investigating single-predictor effects, as this controls the Type I error rate (Forstmeier &  
834 Schielzeth 2011). Reporting the full model also helps reduce publication bias towards  
835 strong effects, providing future meta-analyses with estimates of both significant and  
28

836 non-significant effects (Forstmeier & Schielzeth 2011). Global model reporting should  
837 not replace other model selection methods, but provides a robust measure of how likely  
838 significant effects are to arise by sampling variation alone.

839

## 840 **Practical Issues with Applying Information Theory to Biological Data**

841

### 842 *1. Using All-Subsets Selection*

843 All-Subsets selection is the act of fitting a global model, often containing every possible  
844 interaction, and then fitting every possible nested model. On the surface, all-subsets  
845 might appear to be a convenient and fast way of ‘uncovering’ the causal relationships in  
846 the data. All-subsets selection of enormous global models containing large numbers of  
847 predictors and their interactions makes analyses extremely prone to Type I errors and  
848 ‘overfitted’ models. Burnham & Anderson (2002) caution strongly against all-subsets  
849 selection, and instead advocate ‘hard thinking’ about the hypotheses underlying the  
850 data. If adopting an all subsets approach, it is worth noting the number of models to  
851 consider increases exponentially with the number of predictors, where 5 predictors  
852 require  $2^5$  (32) models to be fitted, whilst 10 predictors requires 1024 models, both  
853 *without* including any interactions.

854 The inflation of Type I error rate through all-subsets selection is simple to  
855 demonstrate. Figure 4 shows the results of a simulation exercise where we created  
856 datasets containing various numbers of continuous and categorical variables, fitted a  
857 global model containing all predictors as main effects and no interactions; and then  
858 performed ASS on that model in the ‘MuMIn’ package in *R*. Note that MuMIn’ refers to  
859 ASS as ‘dredging’ (the ‘dredge’ command), and this *model* dredging is separate from  
860 *data* dredging sensu Burnham & Anderson (2002). All simulated predictors were  
861 samples drawn from populations representing the null hypothesis, i.e. having zero  
862 influence on the response variable. We considered all models with an AIC score of  
863 within 6 of the best-supported AIC model to be equally well supported (also referred to  
864 as the  $\Delta 6$  AIC top model set, Richards 2008) (detailed methods available in Online  
865 Supplementary Material). We assumed a Type I error had occurred when the 95%

29

866 confidence intervals for model averaged parameter estimates from the  $\Delta 6AIC$  set did  
867 not cross zero. The higher the number of terms in the model, the higher the Type I error  
868 rate, reaching a maximum of over 60% probability of falsely including a predictor in the  
869 top model set that was unrelated to the response variable. Importantly, we found that  
870 the rate of increase (slope) in Type I error with added continuous predictors was  
871 modified by the number of categorical variables (Fig. 4), meaning the change in Type 1  
872 error rate per continuous predictor was highest with smaller numbers of categorical  
873 variables. Note that many factors contribute to this high Type I error rate observed here.  
874 For example, just because one level of a factor has 95% intervals that do not span zero  
875 does not mean that the factor as a whole has any explanatory power. See also  
876 Forstmeier & Schielzeth (2011) for a discussion of cryptic testing of multiple hypotheses  
877 in a single model.

878         These results help to illustrate why dredging should not be used, and why global  
879 models should not contain huge numbers of variables and interactions without prior  
880 thought about what the models represent for a study system. In cases where all-subsets  
881 selection from a global model is performed, it is important to view these model selection  
882 exercises as exploratory (Symonds & Moussali 2011), and hold some data back from  
883 these exploratory analyses to be used for cross-validation with the top model(s) (see  
884 Dochterman and Jenkins 2011 and references therein). Here, 90% of the data can be  
885 used to fit the model(s), with the remaining 10% used for confirmatory analysis to  
886 quantify how well the model(s) perform for prediction (Zuur & Ieno 2016). Such an  
887 approach requires a huge amount of data (Dochterman and Jenkins 2011), but cross-  
888 validation to validate a model's predictive ability is rare and should result in more robust  
889 inference (see also Fieberg & Johnson 2015).

890         Therefore, best practice is to consider only a handful of hypotheses and then build a  
891 single statistical model to reflect each hypothesis. This makes inference easier because  
892 the resulting top model set will likely contain fewer parameters, and certainly fewer  
893 spuriously 'significant' parameters (Burnham & Anderson 2002; Arnold 2010). However,  
894 we argue all subsets selection may be sensible in a limited number of circumstances  
895 when testing causal relationships between explanatory variables and the response  
896 variable. For example, if the most complex model contains two main effects and their  
30

897 interaction, performing all subsets selection on that model is identical to building the four  
898 competing models (including the null model) nested in the global model, all of which  
899 may be considered likely to be supported by the data. It is worth remembering that the  
900 Type I error rate can quickly exceed the nominal 5% threshold if these conditions are  
901 not met (Fig. 4). Moreover, a small number of models built to reflect well-reasoned  
902 hypotheses are only valid if the predictors therein are not collinear (see 'Collinearity'  
903 section below). All-subsets selection using the R package *MuMIn* (Bartoń 2016) will not  
904 automatically check for collinearity, and so the onus falls on the researcher to be  
905 thorough in checking for such problems.

906

## 907 2. *Deciding Which Information Criterion To Use*

908 Several information criteria are available to rank competing models, but their  
909 calculations differ subtly. Commonly applied criteria include Akaike's Information  
910 Criterion (AIC), the small sample size correction of AIC for when  $n/k < 40$  (AICc), and the  
911 Bayesian Information Criterion (BIC). QAIC is an adjustment to AIC that accounts for  
912 overdispersion, and should be used when overdispersion has been identified in a model  
913 (see 'Overdispersion section' above). Note QAIC is not required if the overdispersion in  
914 the dataset has been modelled using zero-inflated models, observation-level random  
915 effects, or compound probability distributions. Bolker et al (2009) and Grueber et al  
916 (2011) provide details of how to calculate these criteria.

917 AIC maximises the fit/complexity trade-off of a model by balancing the model fit  
918 with the number of estimated parameters. AICc and BIC both penalise the IC score  
919 based on total sample size  $n$ , but the degree of penalty for AICc is less severe than BIC  
920 for moderate sample sizes, and more severe for very low sample size (Brewer et al  
921 2016). Whilst AIC tend to select overly complex models, Burnham and Anderson (2002)  
922 criticised BIC for selecting overly simplistic models (underfitting). BIC is also criticised  
923 because it operates on the assumption that the true model is in the model set under  
924 consideration, whereas in ecological studies this is unlikely to be true (Burnham &  
925 Anderson 2002; 2004). Issues exist with both AIC and BIC in a GLMM context for  
926 estimating the number of parameters for a random effect (Bolker et al 2009; Grueber et  
927 al 2011), and although degrees of freedom corrections to remedy this problem exist it is  
31



928 not always clear what method is being employed by software packages (see Bolker et al  
929 2009 Box 3). Brewer et al (2016) show how the optimality of AIC, AICc and BIC for  
930 prediction changes with both sample size and effect size of predictors (see also  
931 Burnham and Anderson 2004). Therefore, the choice between the two metrics is not  
932 straightforward, and may depend on the goal of the study i.e. model selection vs  
933 prediction, see Grueber et al 2011 Box 1.

934

### 935 *3. Choice of $\Delta$ AIC Threshold*

936 Once all models have been ranked by an information criterion, it is common practice to  
937 identify a “top model set” containing all models assumed to have comparable support in  
938 the data, normally based on the change in AIC values relative to the best AIC model  
939 ( $\Delta$ AIC). Historically, Burnham & Anderson (2002) recommended that only models with  
940  $\Delta$ AIC between 0-2 should be used for inference, but subsequent work has shown that at  
941 least  $\Delta$ 6 AIC is required to guarantee a 95% probability that the best (expected)  
942 Kullback-Leibler Distance model is in the top model set (Richards 2008; see also  
943 Burnham et al 2011). Alternatively, models can be ranked by their Akaike weights and  
944 all those with an Akaike weight  $\geq 0.95$  retained in the “95% confidence set” (Burnham &  
945 Anderson 2002; Symonds & Moussali 2011). Using high cut-offs is not encouraged, to  
946 avoid overly complex model sets followed by invalid results (Richards 2008; Grueber et  
947 al. 2011) but deciding on how many is too many remains a contentious issue (Grueber  
948 et al. 2011). We suggest  $\Delta$ 6 as a minimum following Richards (2005; 2008).

949

### 950 *4. Using the Nesting Rule to Improve Inference from the Top Model Set*

951 It is well known that AIC tends towards overly complex models (‘overfitting’, Burnham &  
952 Anderson 2002). As AIC only adds a 2 point penalty to a model for inclusion of a new  
953 term, Arnold (2010) demonstrated that adding a nuisance predictor to a well-fitting  
954 model leads to a  $\Delta$ AIC value of the new model of  $\sim 2$ , therefore appearing to warrant  
955 inclusion in the top model set (see section above). Therefore, inference can be greatly  
956 improved by eliminating models from the top model set that are more complex versions  
957 of nested models with better AIC support, known as the nesting rule (Richards 2005;  
958 2008; Richards et al 2011). Doing so greatly reduces the number of models to be used

959 for inference, and improves parameter accuracy (Arnold 2010; Richards et al 2008).  
960 Symonds & Moussali (2011) caution that its applicability has not yet been widely  
961 assessed over a range of circumstances, but the theory behind its application is sound  
962 and intuitive (Arnold 2010). One potential problem is that once models have removed  
963 from the top model set, interpretation of the Akaike weights for the remaining models  
964 becomes difficult, and thus model-averaged estimates using these weights may not be  
965 sensible.

966

#### 967 *5. Using Akaike Weights to Quantify Variable Importance*

968 With a top model set in hand, it is common practice to use the summed Akaike weights  
969 of every model in that set in which a predictor of interest occurs as a measure of  
970 'variable importance' (e.g. Grueber et al 2011). Recent work has demonstrated that this  
971 approach is flawed because Akaike weights are interpreted as relative model  
972 probabilities, and give no information about the importance of individual predictors in a  
973 model (Cade 2015), and fail to distinguish between variables with weak or strong effects  
974 (Galipaud et al 2014; 2017). The sum of Akaike weights as a measure of variable  
975 importance may at best be a measure of how likely a variable would be included after  
976 repeated sampling of the data (Burnham & Anderson 2002; Cade 2015, but see  
977 Galipaud et al 2017). A better measure of variable importance would be to compare  
978 standardised effect sizes (Schielzeth 2010; Cade 2015).

979

#### 980 *6. Model Averaging when Predictors Are Collinear*

981 The aim of model averaging is to incorporate the uncertainty of the size and presence of  
982 effects among a set of candidate models with equal support in the data. Model  
983 averaging using Akaike weights proceeds on the assumption that predictors are on  
984 common scales across models and are therefore comparable. Unfortunately, the nature  
985 of multiple regression means that the scale and sign of coefficients will change across  
986 models depending on the presence or absence of other variables in a focal model  
987 (Cade 2015). The issue of predictor scaling changing across models is particularly  
988 exacerbated when predictors are collinear, even when VIF values are low (Burnham  
989 and Anderson 2002; Lukacs, Burnham & Anderson 2010; Cade 2015). Cade (2015)

33

990 recommends standardising model parameters based on partial standard deviations to  
991 ensure predictors are on common scales across models prior to model averaging  
992 (details in Cade 2015). We stress again the need to assess multicollinearity among  
993 predictors in multiple regression modelling before fitting models (Zuur & Ieno 2016) and  
994 before model-averaging coefficients from those models (Lukacs, Burnham & Anderson  
995 2010; Cade 2015)

996

997

## 998 Conclusion

999 We hope this article will act as both a guide, and as a gateway to further reading, for  
1000 both new researchers and those wishing to update their portfolio of analytic techniques.  
1001 Here we distill our message into a bulleted list.

1002 1. Modern mixed effect models offer an unprecedented opportunity to explore complex  
1003 biological problems by explicitly modelling non-Normal data structures and/or non-  
1004 independence among observational unit. However, the LMM and GLMM toolset should  
1005 be used with caution.

1006 2. Rigorous testing of both model fit ( $R^2$ ) and model adequacy (violation of assumptions  
1007 like homogeneity of variance) must be carried out. We must recognise that satisfactory  
1008 fit does not guarantee we have not violated the assumptions of LMM, and vice versa.  
1009 Interpret measures of  $R^2$  for (G)LMMs with hierarchical errors cautiously, especially  
1010 when OLRE are used.

1011 3. Collinearity among predictors is difficult to deal with and can severely impair model  
1012 accuracy. Be especially vigilant if data are from field surveys rather than controlled  
1013 experiments, as collinearity is likely to be present.

1014 4. Data dredging or 'fishing expeditions' are very risky and inflate the number of false  
1015 positives enormously. Including all combinations of predictors in a model requires strong  
1016 *a priori* justification.

1017 5. When including a large number of predictors is necessary, backwards selection and  
1018 NHST should be avoided, and ranking via AIC of all competing models is preferred. A

1019 critical question that remains to be addressed is whether model selection based on  
1020 information theory is superior to NHST even in cases of balanced experimental designs  
1021 with few predictors.

1022 6. Data simulation is a powerful but underused tool. If the analyst harbours any  
1023 uncertainty regarding the fit or adequacy of the model structure, then the analysis of  
1024 data simulated to recreate the perceived structure of the favoured model can provide  
1025 reassurance, or justify doubt.

1026 7. Wherever possible, provide diagnostic assessment of model adequacy, and metrics  
1027 of model fit, even if in Supplementary Material.

1028 8. Other modelling approaches such as Bayesian inference are available, and allow  
1029 much greater flexibility in choice of model structure, error structure and link function.  
1030 However, the ability to compare among competing models is underdeveloped, and  
1031 where these tools do exist, they are not yet accessible enough to non-experts to be  
1032 useful.

1033

1034

1035

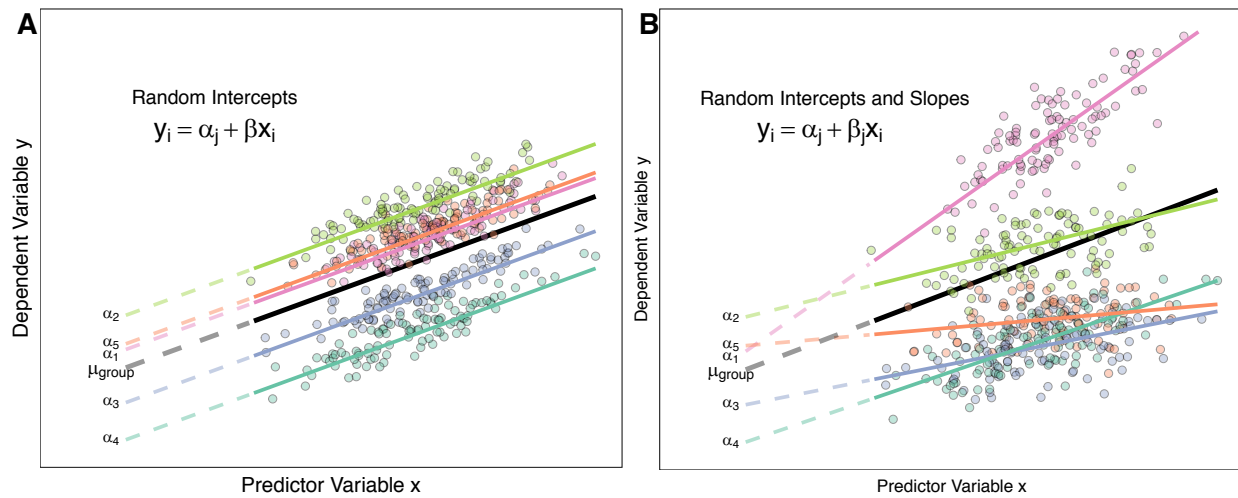
## 1036 Acknowledgements

1037 This paper is the result of a University of Exeter workshop on best practice for the  
1038 application of mixed effects models and model selection in ecological studies.

1039

1040

1041



1042

1043

### Figure 1. Differences between Random Intercept vs Random Slope Models

1044

(A) A random-intercepts model where the outcome variable  $y$  is a function of

1045

predictor  $x$ , with a random intercept for group ID (coloured lines). Because all groups

1046

have been constrained to have a common slope, their regression lines are parallel.

1047

Solid lines are the regression lines fitted to the data. Dashed lines trace the regression

1048

lines back to the  $y$  intercept (0 in this case). Point colour corresponds to group ID of the

1049

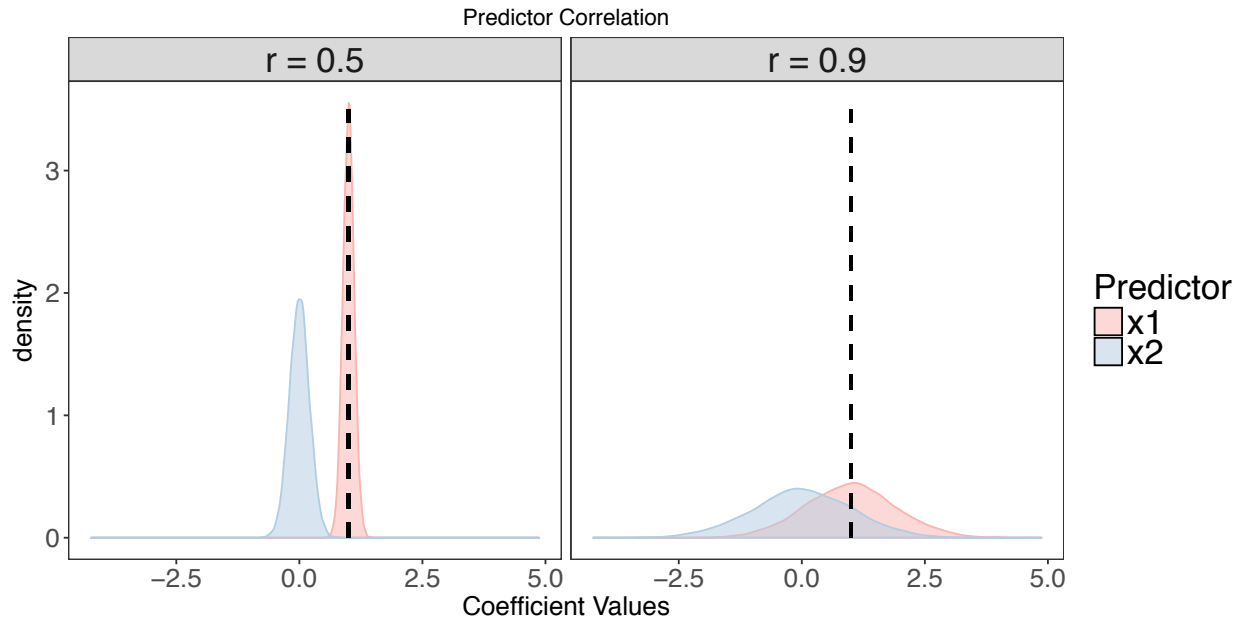
data point. The black line represents the global mean value of the distribution of random

1050

effects.

1051

1052



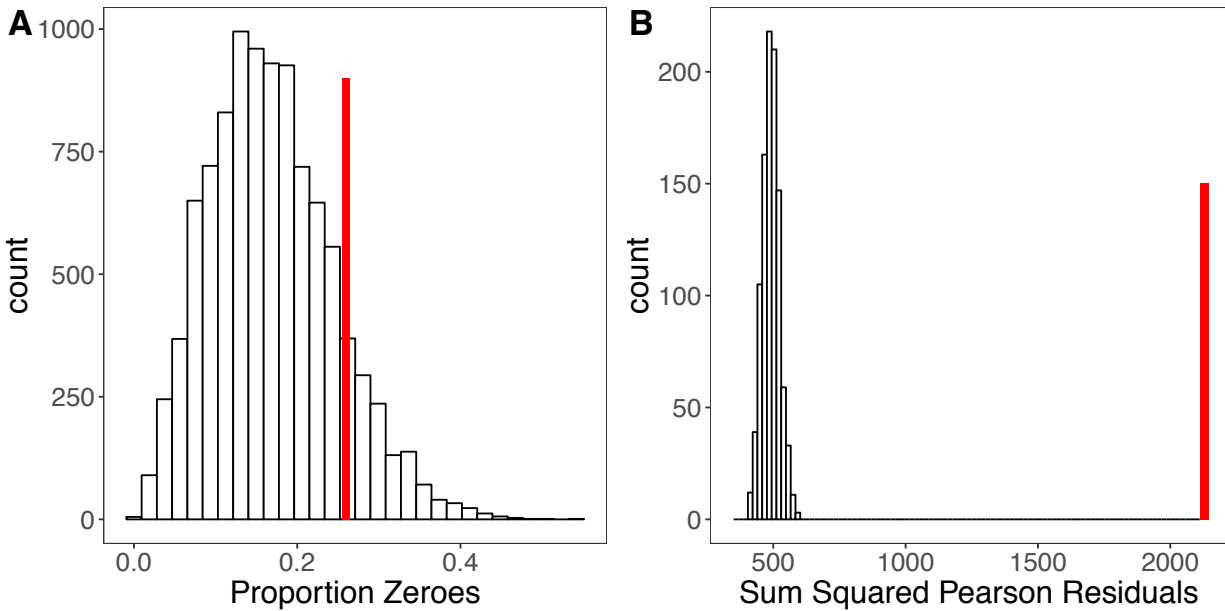
1053

1054 **Figure 2. The effect of collinearity on model parameter estimates.**

1055 We simulated 10,000 iterations of a model  $y \sim x1 + x2$ , where  $x1$  had a positive effect  
1056 on  $y$  ( $\beta_{x1} = 1$ , vertical dashed line).  $x2$  is collinear with  $x1$  with either a moderate ( $r = 0.5$ ,  
1057 A) or strong correlation ( $r = 0.9$ , B). With moderate collinearity, bias in estimation of  
1058  $\beta_{x1}$  is minimal, but variance in estimation of  $\beta_{x2}$  is large. When collinearity is strong, bias  
1059 in estimation of  $\beta_{x1}$  is large, with 14% of simulations estimating a negative coefficient for  
1060 the effect of  $x1$ . For more elaborate versions of these simulations, see Freckleton  
1061 (2011)

1062

1063



1064

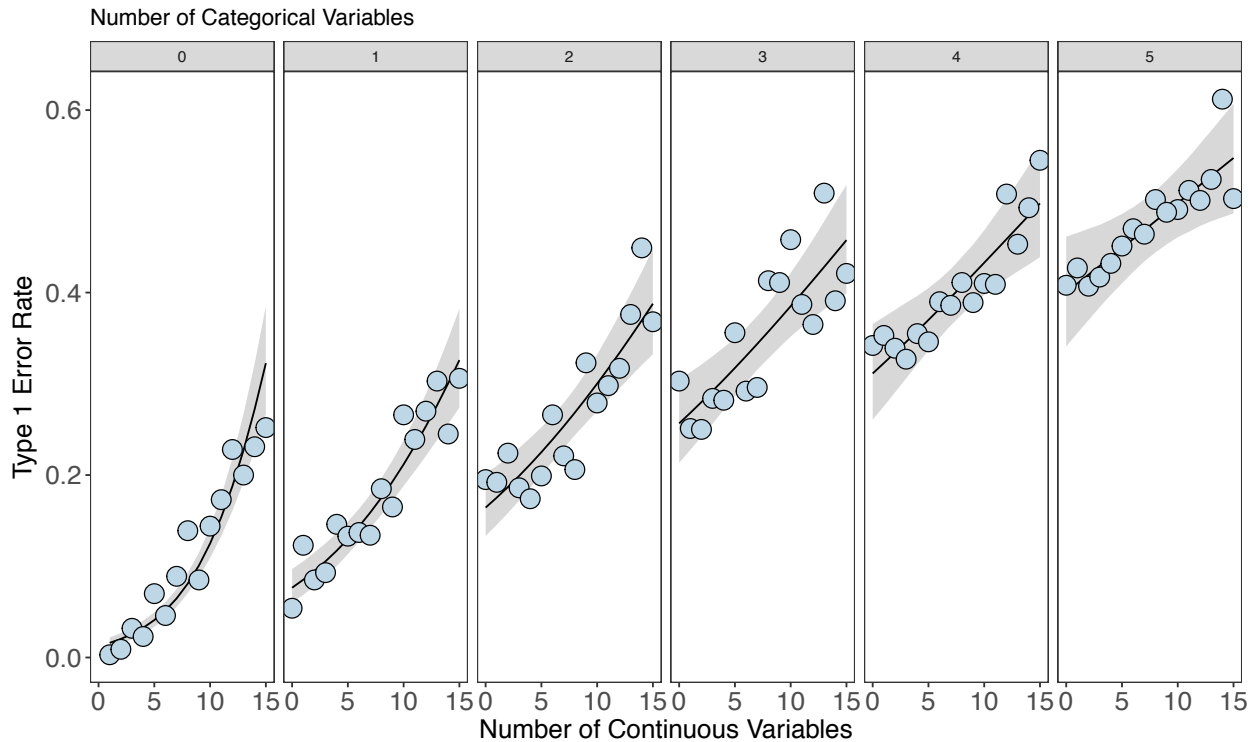
### 1065 **Figure 3. Using Simulation to Assess Model Fit for GLMMs**

1066 (A) Histogram of the proportion of zeroes in 10,000 datasets simulated from a Poisson  
 1067 GLMM. Vertical red line shows the proportion of zeroes in our real dataset. There is no  
 1068 strong evidence of zero-inflation for these data. (B) Histogram of the sum of squared  
 1069 Pearson residuals for 1000 parametric bootstraps where the Poisson GLMM has been  
 1070 re-fitted to the data at each step. Vertical red line shows the test statistic for the original  
 1071 model, which lies well outside the simulated frequency distribution. The ratio of the real  
 1072 statistic to the simulated data can be used to calculate a mean dispersion statistic and  
 1073 95% confidence intervals, which for these data is mean 3.16, 95% CI 2.77 – 3.59.  
 1074 Simulating from models provides a simple yet powerful set of tools for assessing model  
 1075 fit and robustness.

1076

1077

1078



1079

1080 **Figure 4. The effect of data dredging on Type 1 Error Rate as a function of the**  
1081 **number of continuous and categorical variables included in the global model**

1082 Adding both categorical and continuous predictors to the models (increasing complexity)  
1083 increases the Type I error rate (95% confidence intervals of model averaged parameter  
1084 estimates do not cross zero). The slope of the increase in Type I error rate with increase  
1085 in the number of continuous predictors is modified by how many categorical predictors  
1086 there are in the model, with steeper increases in Type 1 error rate for lower numbers of  
1087 categorical predictors. However, the Type I error rate was highest overall for global  
1088 models containing the largest numbers of parameters. For full details of the simulation  
1089 methodology, see supplementary file S1).

1090

1091



1092

1093 **References**

- 1094 Allegue H, Araya-Ajoy YG, Dingemanse NJ, Dochtermann NA, Garamszegi LZ,  
1095 Nakagawa S, Reale D, Schielzeth H, Westneat DF. 2017. Statistical Quantification  
1096 of Individual Differences (SQUID): an educational and statistical tool for  
1097 understanding multilevel phenotypic data in linear mixed models. *Methods in*  
1098 *Ecology and Evolution* 8:257-67.
- 1099 Arnold TW. 2010. Uninformative parameters and model selection using Akaike's  
1100 Information Criterion. *The Journal of Wildlife Management* 74: 1175-1178.
- 1101 Austin MP. 2002. Spatial prediction of species distribution: an interface between  
1102 ecological theory and statistical modelling. *Ecological Modelling* 157: 101–118.
- 1103 Barker RJ, Link WA. 2015. Truth, models, model sets, AIC, and multimodel inference: A  
1104 Bayesian perspective. *The Journal of Wildlife Management* 79: 730–738.
- 1105 Barr DJ, Levy R, Scheepers C, Tily HJ. 2013. Random effects structure for confirmatory  
1106 hypothesis testing: Keep it maximal. *Journal of memory and language* 68:255-78.
- 1107 Bartoń K. 2016. MuMIn: Multi-Model Inference. R package version  
1108 1.15.6.<https://CRAN.R-project.org/package=MuMIn>
- 1109 Bates D, Maechler M, Bolker B, Walker S. 2015. Fitting Linear Mixed-Effects Models  
1110 Using lme4. *Journal of Statistical Software* 67: 1-48.
- 1111 Bates D, Kliegl R, Vasishth S, Baayen H. 2015. Parsimonious mixed models. *arXiv*  
1112 *preprint arXiv:1506.04967*.
- 1113 Bolker BM, Brooks ME, Clark CJ, Geange SW, Poulsen JR, Stevens MHH, White JSS.  
1114 2009. Generalized linear mixed models: a practical guide for ecology and  
1115 evolution. *Trends in Ecology and Evolution* 24: 127–135.
- 1116 Breslow NE, Clayton DG. 1993. Approximate inference in generalized linear mixed  
1117 models. *Journal of the American statistical Association* 88: 9-25.
- 1118 Brewer MJ, Butler A, Cooksley SL. 2016. The relative performance of AIC, AICC and  
1119 BIC in the presence of unobserved heterogeneity. *Methods in Ecology and*  
1120 *Evolution* 7: 679-692.

40

- 1121 Burnham KP, Anderson DR. 2002. Model Selection and Multimodel Inference: A  
1122 Practical Information-Theoretic Approach, Second. Springer-Verlag, New York.
- 1123 Burnham KP, Anderson DR. 2004. Multimodel inference: understanding AIC and BIC in  
1124 model selection. *Sociological Methods & Research* 33: 261-304.
- 1125 Burnham KP, Anderson DR, Huyvaert KP. 2011. AIC model selection and multimodel  
1126 inference in behavioral ecology: Some background, observations, and  
1127 comparisons. *Behavioral Ecology and Sociobiology* 65: 23–35.
- 1128 Cade BS. 2015. Model averaging and muddled multimodel inferences. *Ecology* 96:  
1129 2370–2382.
- 1130 Chatfield C. 1995. Model uncertainty, data mining and statistical inference (with  
1131 discussion). *Journal of the Royal Statistical Society, Series A* 158: 419-66.
- 1132 Cox DR, Snell EJ. 1989. *The Analysis of Binary Data*, 2nd ed. London: Chapman and  
1133 Hall.
- 1134 Crawley (2013) *The R Book*. Second Edition. Wiley, Chichester UK.
- 1135 Dochtermann NA, Jenkins SH. 2011. Developing multiple hypotheses in behavioural  
1136 ecology. *Behavioral Ecology and Sociobiology* 65: 37-45.
- 1137 Dominicus A, Skrondal A, Gjessing HK, Pedersen NL, Palmgren J. 2006. Likelihood ratio  
1138 tests in behavioral genetics: problems and solutions. *Behavior Genetics* 36: 331–  
1139 340.
- 1140 Dormann CF, Elith J, Bacher S, Buchmann C, Carl G, Carré G, Marquéz JR, Gruber B,  
1141 Lafourcade B, Leitão PJ, Münkemüller T. 2013. Collinearity: a review of methods  
1142 to deal with it and a simulation study evaluating their performance. *Ecography* 36:  
1143 027–046.
- 1144 Ellison AM. 2004. Bayesian inference in ecology. *Ecology letters* 7: 509-520.
- 1145 Elston, DA, Moss R, Boulinier T, Arrowsmith C, Lambin X, 2001. Analysis of  
1146 aggregation, a worked example: numbers of ticks on red grouse  
1147 chicks. *Parasitology* 122: 563-569.
- 1148 Fieberg J, Johnson DH. 2015. MMI: Multimodel inference or models with management  
1149 implications? *The Journal of Wildlife Management* 79: 708–718.

- 1150 Forstmeier W, Schielzeth H. 2011. Cryptic multiple hypotheses testing in linear models:  
1151 Overestimated effect sizes and the winner's curse. *Behavioral Ecology and*  
1152 *Sociobiology* 65: 47–55.
- 1153 Freckleton RP. 2011. Dealing with collinearity in behavioural and ecological data: model  
1154 averaging and the problems of measurement error. *Behavioral Ecology and*  
1155 *Sociobiology* 65: 91-101.
- 1156 Galipaud M, Gillingham MAF, David M, Dechaume-Moncharmont FX. 2014. Ecologists  
1157 overestimate the importance of predictor variables in model averaging: a plea for  
1158 cautious interpretations. *Methods in Ecology and Evolution* 5, 983-991.
- 1159 Galipaud M, Gillingham MAF, Dechaume-Moncharmont FX. 2017. A farewell to the sum  
1160 of Akaike weights: The benefits of alternative metrics for variable importance  
1161 estimations in model selection. *Methods in Ecology and Evolution* 00:1–11.  
1162 <https://doi.org/10.1111/2041-210X.12835>
- 1163 Gelman A, Hill J. 2007. Data analysis using regression and hierarchical/multilevel  
1164 models. New York, NY, USA: Cambridge University Press.
- 1165 Gelman A. 2008. Scaling regression inputs by dividing by two standard  
1166 deviations. *Statistics in Medicine* 27: 2865-2873.
- 1167 Gelman A, Pardoe I. 2006. Bayesian measures of explained variance and pooling in  
1168 multilevel (hierarchical) models. *Technometrics* 48: 241-251.
- 1169 Graham ME (2003) Confronting multicollinearity in multiple linear regression. *Ecology*  
1170 84: 2809-2815
- 1171 Grueber CE, Nakagawa S, Laws RJ, Jamieson IG. 2011. Multimodel inference in  
1172 ecology and evolution: Challenges and solutions. *Journal of Evolutionary Biology*  
1173 24: 699–711.
- 1174 Harrison XA. 2014. Using observation-level random effects to model overdispersion in  
1175 count data in ecology and evolution. *PeerJ* 2: e616.
- 1176 Harrison XA. 2015. A comparison of observation-level random effect and Beta-Binomial  
1177 models for modelling overdispersion in Binomial data in ecology &  
1178 evolution. *PeerJ*, 3: p.e1114.
- 1179 Halsey LG, Curran-Everett D, Vowler SL, Drummond GB. 2015. The fickle P value  
1180 generates irreproducible results. *Nature Methods* 12: 179-185.

- 1181 Hegyi G, Garamszegi LZ. 2011. Using information theory as a substitute for stepwise  
1182 regression in ecology and behaviour. *Behavioral Ecology and Sociobiology* 65: 69-  
1183 76.
- 1184 Hilbe JM. 2011. *Negative binomial regression*. Cambridge University Press.
- 1185 Houslay T, Wilson A. 2017. Avoiding the misuse of BLUP in behavioral ecology.  
1186 *Behavioral Ecology* arx023 doi:10.1093/beheco/arx023
- 1187 Ives AR. 2015. For testing the significance of regression coefficients, go ahead and log-  
1188 transform count data. *Methods in Ecology and Evolution* 6:, 828-835.
- 1189 James FC, McCullugh CF. 1990. Multivariate Analysis In Ecology And Systematics:  
1190 Panacea Or Pandora Box. *Annual Review of Ecology and Systematics* 21: 129–  
1191 166.
- 1192 Johnson JB, Omland KS. 2004. Model selection in ecology and evolution. *Trends in*  
1193 *Ecology and Evolution* 19: 101–108.
- 1194 Johnson PCD. 2014. Extension of Nakagawa & Schielzeth's  $R^2$  GLMM to random  
1195 slopes models. *Methods in Ecology and Evolution* 5: 944-946.
- 1196 Kass RE, Caffo BS, Davidian M, Meng XL, Yu B, Reid N. 2016. Ten simple rules for  
1197 effective statistical practice. *PLoS computational biology* 12: p.e1004961.
- 1198 Keene ON. 1995. The log transform is special. *Statistics in Medicine* 14: 811–819.
- 1199 Kéry M. 2010. Introduction to WinBUGS for ecologists: Bayesian approach to  
1200 regression, ANOVA, mixed models and related analyses. Academic Press.
- 1201 Kuznetsova A, Brockhoff PB, Christensen RHB. 2014. Package 'lmerTest'. Test for  
1202 random and fixed effects for linear mixed effect models (lmer objects of lme4  
1203 package). R package ver.2.
- 1204 Lefcheck JS. 2015. piecewiseSEM: Piecewise structural equation modeling in R for  
1205 ecology, evolution, and systematics. *Methods in Ecology and Evolution* 7: 573-  
1206 579.
- 1207 Lindberg MS, Schmidt JH, Walker J. 2015. History of multimodel inference via model  
1208 selection in wildlife science. *The Journal of Wildlife Management* 79: 704–707.
- 1209 Low-Décarie E, Chivers C, Granados M. 2014. Rising complexity and falling explanatory  
1210 power in ecology. *Frontiers in Ecology and the Environment* 12: 412-418.

- 1211 Lüdecke D. 2017. SjPlot: Data Visualization for Statistics in Social Science. 2017 R  
1212 *package version, 2.4.0.*
- 1213 Lukacs PM, Burnham KP, Anderson DR. 2010. Model selection bias and Freedman's  
1214 paradox. *Annals of the Institute of Statistical Mathematics* 62: 117–125.
- 1215 Mundry R. 2011. Issues in information theory-based statistical inference—a  
1216 commentary from a frequentist's perspective. *Behavioral Ecology and*  
1217 *Sociobiology* 65: 57-68.
- 1218 Murtaugh PA. 2007. Simplicity and complexity in ecological data analysis. *Ecology* 88:  
1219 56-62.
- 1220 Murtaugh PA. 2009. Performance of several variable-selection methods applied to real  
1221 ecological data. *Ecology Letters* 10: 1061-1068.
- 1222 Murtaugh PA. 2014. In defense of P values. *Ecology* 95: 611-617
- 1223 Nagelkerke NJ. 1991. A note on a general definition of the coefficient of determination.  
1224 *Biometrika* 78: 691-692.
- 1225 Nakagawa S, Foster T. 2004. The case against retrospective statistical power analyses  
1226 with an introduction to power analysis. *Acta Ethologica* 7: 103-108.
- 1227 Nakagawa S, Freckleton RP. 2008. Missing inaction: the dangers of ignoring missing  
1228 data. *Trends in Ecology and Evolution* 23(11): 592-596.
- 1229 Nakagawa S, Freckleton RP. 2011. Model averaging, missing data and multiple  
1230 imputation: a case study for behavioural ecology. *Behavioral Ecology and*  
1231 *Sociobiology* 65: 103-116.
- 1232 Nakagawa S, Schielzeth H. 2010. Repeatability for Gaussian and non-Gaussian data: a  
1233 practical guide for biologists. *Biological Reviews* 85: 935-956
- 1234 Nakagawa S, Schielzeth H. 2013. A general and simple method for obtaining R<sup>2</sup> from  
1235 generalized linear mixed-effects models. *Methods in Ecology and Evolution* 4:  
1236 133-142.
- 1237 Nakagawa S., Johnson PC, Schielzeth H. 2017. The coefficient of determination R<sup>2</sup> and  
1238 intra-class correlation coefficient from generalized linear mixed-effects models  
1239 revisited and expanded. *Journal of The Royal Society Interface* 14(134),  
1240 p.20170213.

- 1241 Nickerson RS. 2000. Null Hypothesis Significance Testing: A Review of an Old and  
1242 Continuing Controversy. *Psychological Methods* 5: 241-301.
- 1243 O'Hara RB, Kotze DJ. 2010. Do not log-transform count data. *Methods in Ecology and*  
1244 *Evolution* 1: 118-122.
- 1245 Peters RH. 1991. *A critique for ecology*. Cambridge University Press.
- 1246 Peig J, Green AJ. 2009. New perspectives for estimating body condition from  
1247 mass/length data: the scaled mass index as an alternative method. *Oikos* 118:  
1248 1883-1891.
- 1249 Quinn GP, Keough MJ. 2002. *Experimental design and data analysis for biologists*.  
1250 Cambridge University Press.
- 1251 R Core Team. 2016. R: A language and environment for statistical computing. R  
1252 Foundation for Statistical Computing, Vienna, Austria. URL [https://www.R-](https://www.R-project.org/)  
1253 [project.org/](https://www.R-project.org/).
- 1254 Richards SA. 2005. Testing ecological theory using the information-theoretic approach:  
1255 examples and cautionary results. *Ecology* 86: 2805-2814.
- 1256 Richards SA. 2008. Dealing with overdispersed count data in applied ecology. *Journal*  
1257 *of Applied Ecology* 45 218–227.
- 1258 Richards, SA, Whittingham MJ, Stephens PA. 2011. Model selection and model  
1259 averaging in behavioural ecology: the utility of the IT-AIC framework. *Behavioral*  
1260 *Ecology and Sociobiology* 65: 77–89.
- 1261 Rykiel EJ. 1996. Testing ecological models: The meaning of validation. *Ecological*  
1262 *Modelling* 90: 229-244.
- 1263 Satterthwaite FE. 1946. An approximate distribution of estimates of variance  
1264 components. *Biometrics Bulletin* 2(6): 110-114.
- 1265 Scheipl F, & Bolker, B. 2016. RLRsim: Exact (Restricted) Likelihood Ratio Tests for  
1266 Mixed and Additive Models *Computational Statistics & Data Analysis*. R package  
1267 version 3.1-3. <https://cran.r-project.org/web/packages/RLRsim/index.html>
- 1268 Schielzeth H, Forstmeier W. 2009. Conclusions beyond support: overconfident  
1269 estimates in mixed models. *Behavioral Ecology* 20: 416-420.

- 1270 Schielzeth H, Nakagawa S. 2013. Nested by design: model fitting and interpretation in a  
1271 mixed model era. *Methods in Ecology Evolution* 4: 14-24
- 1272 Schielzeth H. 2010. Simple means to improve the interpretability of regression  
1273 coefficients. *Methods in Ecology and Evolution* 1: 103-113
- 1274 Southwood TRE, Henderson PA. 2000. *Ecological methods*. John Wiley &  
1275 Sons. Stephens PA, Buskirk SW, Hayward GD, Martinez Del Rio C. 2005.  
1276 Information theory and hypothesis testing: a call for pluralism. *Journal of Applied*  
1277 *Ecology* 42: 4-12.
- 1278 Symonds MRE, Moussalli A. 2011. A brief guide to model selection, multimodel  
1279 inference and model averaging in behavioural ecology using Akaike's information  
1280 criterion. *Behavioral Ecology and Sociobiology* 65: 13–21.
- 1281 Vaida F, Blanchard S. 2005. Conditional Akaike information for mixed-effects models.  
1282 *Biometrika* 92: 351–370
- 1283 van de Pol M, Wright J. 2009. A simple method for distinguishing within-versus  
1284 between-subject effects using mixed models. *Animal Behaviour* 77: 753-758.
- 1285 Verbenke G, Molenberghs G. 2000. Linear mixed models for longitudinal data. New  
1286 York, Springer.
- 1287 Warton D, Hui F. 2011. The arcsine is asinine: the analysis of proportions in ecology.  
1288 *Ecology* 92: 3-10
- 1289 Warton DI, Lyons M, Stoklosa J, Ives AR. 2016. Three points to consider when  
1290 choosing a LM or GLM test for count data. *Methods in Ecology and Evolution* 7:  
1291 882-90.
- 1292 Wilson AJ, Réale D, Clements MN, Morrissey MM, Postma E, Walling CA, Kruuk LEB,  
1293 Nussey DH. 2010. An ecologist's guide to the animal model. *Journal of Animal*  
1294 *Ecology* 79: 13–26.
- 1295 Wood SN, Goude Y, Shaw S. 2015. Generalized additive models for large data  
1296 sets. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 64:139-  
1297 155.

- 1298 Whittingham MJ, Stephens PA, Bradbury RB, Freckleton RP 2006. Why do we still use  
1299 stepwise modelling in ecology and behaviour? *Journal of Animal Ecology* 75:  
1300 1182-1189.
- 1301 Zuur AF, Ieno EN, Walker NJ, Saveliev AA, Smith GM. 2009 *Mixed Effects Models and*  
1302 *Extensions in Ecology with R* Springer, New York
- 1303 Zuur AF, Ieno EN, Elphick CS. 2010. A protocol for data exploration to avoid common  
1304 statistical problems. *Methods in Ecology and Evolution* 1: 3-14.
- 1305 Zuur AF, Ieno EN, 2016. A protocol for conducting and presenting results of regression-  
1306 type analyses. *Methods in Ecology and Evolution* 7: 636-645.
- 1307
- 1308