

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>

Best practice in mixed effects modelling and multi-model inference in ecology

Xavier A Harrison ^{Corresp., 1}, Lynda Donaldson ², Maria Eugenia Correa-Cano ², Julian Evans ^{3,4}, David N Fisher ^{3,5}, Cecily Goodwin ², Beth Robinson ⁶, David J Hodgson ³, Richard Inger ^{2,3}

¹ Institute of Zoology, Zoological Society of London, London, United Kingdom

² Environment and Sustainability Institute, University of Exeter, Penryn, United Kingdom

³ Centre for Ecology and Conservation, University of Exeter, Penryn, United Kingdom

⁴ Department of Biology, University of Ottawa, Ottawa, Canada

⁵ Department of Integrative Biology, University of Guelph, Guelph, Canada

⁶ WildTeam Conservation, Surfside, St Merryn, Padstow PL28 8NU, United Kingdom

Corresponding Author: Xavier A Harrison

Email address: x.harrison@ucl.ac.uk

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.

1 Best Practice in Mixed Effects Modelling and Multi-model Inference in Ecology

2

3 Xavier A. Harrison¹, Lynda Donaldson², Maria Eugenia Correa-Cano², Julian Evans^{3,4},

4 David N. Fisher^{3&5}, Cecily E. D. Goodwin², Beth S. Robinson^{2&6}, Dave Hodgson³ and

5 Richard Inger^{2&3}.

6

7 ¹ Institute of Zoology, Zoological Society of London, London NW1 4RY

8 ² Environment and Sustainability Institute, University of Exeter

9 ³ Centre for Ecology and Conservation, University of Exeter

10 ⁴ Department of Biology, University of Ottawa, Canada

11 ⁵ Department of Integrative Biology, University of Guelph

12 ⁶ WildTeam Conservation, Surfside, Trelantis, St Merryn, Padstow PL28 8NU, Cornwall,

13 UK

14

15

16 Corresponding Author:

17 Xavier Harrison

18

19 Corresponding Email:

20 xav.harrison@gmail.com

21

22

23

24

25

26

27

28

29

30

31

32

33 Introduction

34

35 In recent years, the suite of statistical tools available to biologists and the complexity of
36 biological data analyses have grown in tandem (Low-Decarie et al 2014; Zuur et al
37 2016; Kass et al 2016). The availability of novel and sophisticated statistical techniques
38 means we are better equipped than ever to extract signal from noisy biological data, but
39 it remains challenging to know how to apply these tools, and which statistical
40 technique(s) might be best suited to answering specific questions (Kass et al 2016).
41 Often, simple analyses will be sufficient (Murtaugh 2007), but more complex data
42 structures often require more complex methods such as linear mixed effects models
43 (Zuur et al 2009), generalized additive models (Wood et al 2015) or Bayesian inference
44 (Ellison 2004). Both accurate parameter estimates and robust biological inference
45 require that ecologists be aware of the pitfalls and assumptions that accompany these
46 techniques and adjust modelling decisions accordingly (Bolker et al 2009).

47 Linear mixed effects models (LMMs) and generalized linear mixed effects models
48 (GLMMs), have gained significant traction in the last decade (Zuur et al 2009; Bolker et
49 al 2009). Both extend traditional linear models to include a combination of fixed and
50 random effects as predictor variables. The introduction of random effects affords several
51 non-exclusive benefits. First, biological datasets are often highly structured, containing
52 clusters of non-independent observational units that are hierarchical in nature, and
53 LMMs allow us to explicitly model the non-independence in such data. For example, we
54 might measure several chicks from the same clutch, and several clutches from different
55 females, or we might take repeated measurements of the same chick's growth rate over
56 time. In both cases, we might expect that measurements within a statistical unit (here,
57 an individual, or a female's clutch) might be more similar than measurements from
58 different units. Explicit modelling of the random effects structure will aid correct
59 inference of fixed effects, depending on which level of the system's hierarchy is being
60 manipulated. In our example, if the fixed effect varies or is manipulated at the level of
61 the clutch, then pseudoreplicated measurements of each chick can be controlled

62 carefully using random effects. Alternatively, if fixed effects vary at the level of the chick,
63 then non-independence among clutches or mothers can be accounted for. Random
64 effects typically represent some grouping variable (Breslow and Clayton 1993) and
65 allow the estimation of variance in the response variable within and among these
66 groups. This reduces the probability of false positives (Type I error rates) and false
67 negatives (Type II error rates, e.g. Crawley 2013). Second, inferring the magnitude of
68 variation within and among statistical clusters or hierarchical levels can be highly
69 informative in its own right. In our bird example, understanding whether there is more
70 variation in a focal trait among females within a population, rather than among
71 populations, might be a central goal of the study.

72 LMMs are powerful yet complex tools. Software advances have made these tools
73 accessible to the non-expert and have become relatively straightforward to fit in widely
74 available statistical packages such as R (R Core Team 2016). However, despite this
75 ease of implementation, the correct use of LMMs in the biological sciences is
76 challenging for several reasons: i) they make additional assumptions about the data to
77 those made in more standard statistical techniques such as general linear models
78 (GLMs), and these assumptions are often violated (Bolker et al 2009); ii) interpreting
79 model output correctly can be challenging, especially for the variance components of
80 random effects (Bolker et al 2009; Zuur et al 2009); iii) model selection for LMMs
81 presents a unique challenge, irrespective of model selection philosophy, because of
82 biases in the performance of some tests (e.g. Wald tests, AIC comparisons) introduced
83 by the presence of random effects (Vaida & Blanchard 2005; Dominicus et al 2006;
84 Bolker et al 2009). Collectively, these issues mean that the application of LMM
85 techniques to biological problems can be risky and difficult for those that are unfamiliar
86 with them. There have been several excellent papers in recent years on the use of
87 generalized linear mixed effects models (GLMMs) in biology (Bolker et al 2009), the use
88 of information theory and multi-model inference for studies involving LMMs (Grueber et
89 al 2011), best practice for data exploration (Zuur et al 2009) and for conducting
90 statistical analyses for complex datasets (Zuur et al 2016; Kass et al 2016). At the
91 interface of these excellent guides lies the theme of this paper: an updated guide for the
92 uninitiated through the model fitting and model selection processes when using LMMs.

93 A secondary but no less important aim of the paper is to bring together several key
94 references on the topic of LMMs, and in doing so act as a portal into the primary
95 literature that derives, describes and explains the complex modelling elements in more
96 detail.

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

108 Understanding Fixed and Random Effects

109
110 A key decision of the modelling process is specifying model predictors as fixed or
111 random effects. Unfortunately, the distinction between the two is not always obvious,
112 and is not helped by the presence of multiple, often confusing definitions in the literature
113 (see Gelman and Hill 2007 p. 245). Absolute rules for how to classify something as a
114 fixed or random effect generally are not useful because that decision can change
115 depending on the goals of the analysis (Gelman and Hill 2007). We can illustrate the
116 difference between fitting something as a fixed (M1) or a random effect (M2) using a
117 simple example of a researcher who takes measurements of mass from 100 animals
118 from each of 5 different groups ($n=500$) with a goal of understanding differences among
119 groups in mean mass. We use notation equivalent to fitting the proposed models in the
120 statistical software *R* (R Core Team 2016), with the LMMs fitted using the R package
121 *lme4* (Bates et al. 2015):

122

```
123           M1 <- glm (mass ~ group)
124           M2 <- lmer(mass ~ 1 + (1|group))
```

125

126 Fitting 'group' as a fixed effect in model M1 assumes the 5 'group' means are all
127 independent of one another, and share a common residual variance. Conversely, fitting
128 group as a random intercept model in model M2 assumes that the 5 measured group
129 means are only a subset of the realised possibilities drawn from a 'global' set of
130 population means that follow a Normal distribution with its own mean (μ_{group} , Fig. 1A)
131 and variance (σ^2_{group}). Therefore, LMMs model the variance hierarchically, estimating
132 first the process generating among-group variation in means, and subsequently
133 variation within groups. Treating groups from a field survey as only a subset of the
134 *possible* groups that could be sampled is quite intuitive, because there are likely many
135 more groups (e.g. populations) of the study species in nature than the 5 the researcher
136 measured. Conversely if one has designed an experiment to test the effect of three
137 different temperature regimes on growth rate of plants, specifying temperature
138 treatment as a fixed effect appears sensible because experimenter has deliberately set
139 the variable at a given value of interest. That is, there are no unmeasured groups with
140 respect to that particular experimental design.

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

145

146 *Controlling for non-independence among data points*

147 This is one of the most common uses of a random effect. Complex biological data sets
148 often contain nested and/or hierarchical structures such as repeat measurements from
149 individuals within and across units of time. Random effects allow you to control for this
150 non-independence by constraining non-independent 'units' to have the same intercept
151 and/or slope (Zuur et al 2009; Zuur et al 2016). Whether you fit only random intercepts
152 or both random intercepts and slopes will be decided by the goals of the analysis, and
153 the dependency structure of the data (Zuur et al 2016). Fitting *only* a random intercept

154 means you allow group means to vary, but assume all groups have a common slope for
155 a fitted covariate (fixed effect). Fitting random intercepts *and* slopes means you allow
156 the slope of a predictor to vary based on a separate grouping variable. For example,
157 one hypothesis might be that the probability of successful breeding for an animal is a
158 function of its body mass. If we had measured animals from multiple sampling sites, we
159 might wish to fit 'sampling site' as a random intercept, and estimate a common slope
160 (change in breeding success) for body mass across all sampling sites by fitting it as a
161 fixed effect:

162

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

164

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

170

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

173

174 Schielzeth & Forstmeier (2009) warn that constraining groups to share a common slope
175 can inflate Type I and Type II errors. Consequently, Grueber et al (2011) recommend
176 always fitting both random slopes and intercepts where possible. Whether this is
177 feasible or not will depend on your data structure (see 'Costs to Fitting Random Effects'
178 section below). Figure 1 describes the differences between random intercept models
179 and those also containing random slopes.

180 *Further reading: Zuur et al (2016) shows examples of the difficulties in identifying*
181 *the dependency structure of data and how to use flow charts / graphics to help decide*
182 *model structure. Kery (2010, Ch 12) has an excellent demonstration of how to fit*
183 *random slopes, and how model assumptions change depending on whether you specify*
184 *a correlation between random slopes and intercepts or not. Schielzeth & Forstmeier*

185 (2009) and van de Pol & Wright (2009) are useful references for understanding the
186 utility of random slope models.

187

188 *To improve accuracy of parameter estimation*

189 Random effect models use data from all the groups to estimate the mean and variance
190 of the global distribution of group means. Assuming all group means are drawn from a
191 common distribution causes the estimates of their means to drift towards the global
192 mean μ_{group} . This phenomenon, known as *shrinkage* (Gelman & Hill 2007; Kery 2010),
193 can also lead to smaller and more precise standard errors around means. Shrinkage is
194 strongest for groups with small sample sizes, as the paucity of within-group information
195 to estimate the mean is counteracted by the model using data from other groups to
196 improve the precision of the estimate. This 'partial pooling' of the estimates is a principal
197 benefit of fitting something as a random effect (Gelman & Hill 2007). However, it can
198 feel strange that group means should be shrunk towards the global mean, especially for
199 researchers more used to treating sample means as independent fixed effects.
200 Accordingly, one issue is that variance estimates can be hugely imprecise when there
201 are fewer than 5 levels of the random grouping variable (intercept or slope; see Harrison
202 2015). However, thanks to the Central Limit Theorem, the assumption of Gaussian
203 distribution of group means is usually a good one, and the benefits of hierarchical
204 analysis will outweigh the apparent costs of shrinkage.

205

206 *To estimate variance components*

207 In some cases, the variation among groups will be of interest to ecologists. For
208 example, imagine we had measured the clutch masses of 30 individual birds, each of
209 which had produced 5 clutches ($n=150$). We might be interested in asking whether
210 different females tend to produce consistently different clutch masses (high among-
211 female variance for clutch mass). To do so, we might fit the following model with Clutch
212 Mass as the response variable, no fixed effects, and a Gaussian error structure:

213

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

215

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

233

234 *To make predictions for unmeasured groups*

235 Fixed effect estimates prevent us from making predictions for new groups because the
236 model estimates are only relevant to groups in our dataset (e.g. Zuur et al 2009 p. 327).
237 Conversely, we can use the estimate of the global distribution of population means to
238 predict for the average group using the mean of the distribution μ_{group} for a random
239 effects model (see Fig. 1). We could also sample hypothetical groups from our random
240 effect distribution, as we know its mean and SD (Zuur et al 2016). Therefore, whether
241 something is fitted as a fixed or random effect can depend on the goal of the analysis:
242 are we only interested in the mean values for each group in our dataset, or do we wish
243 to use our results to extend our predictions to new groups? Even if we do not want to
244 predict to new groups, we might wish to fit something as a random effect to take
245 advantage of the shrinkage effect and improved parameter estimation accuracy.

246

247 **Considerations When Fitting Random Effects**

248 Random effect models have several desirable properties (see above), but their use
249 comes with some caveats. First, they are quite ‘data hungry’; as a rule, you need at
250 least 5 ‘levels’ (groups) for a random intercept term to achieve robust estimates of
251 variance (Gelman & Hill 2007; Harrison 2015). With <5 levels, the mixed model may not
252 be able to estimate the among-population variance accurately. In this case, the variance
253 estimate will either collapse to zero, making your model equivalent to an ordinary GLM
254 (Gelman & Hill 2007 p. 275) or be non-zero but incorrect if the small number of groups
255 you have sampled are not representative of true distribution of means (Harrison 2015).
256 Second, models can be unstable if sample sizes across groups are highly unbalanced
257 i.e. if some groups contain very few data. These issues are especially relevant to
258 random slope models (Grueber et al 2011). Third, an important issue is the difficulty in
259 deciding the “significance” or “importance” of variance among groups. The variance of a
260 random effect is inevitably at least zero, but how big does it need to be to be considered
261 of interest? Fitting a factor as a fixed effect provides a statement of the significance of
262 differences (variation) among groups relatively easily. Testing differences among levels
263 of a random effect is made much more difficult for frequentist analyses, though not so in
264 a Bayesian framework (Kery 2010, see ‘*Testing Significance of Random Effects*’
265 section). Finally, an issue that is not often addressed is that of mis-specification of
266 random effects. GLMMs are powerful tools, but incorrectly parameterising the random
267 effects in your model could yield model estimates that are as unreliable as ignoring the
268 need for random effects altogether. An example would be failure to recognise
269 nonindependence caused by nested structures in the data e.g. multiple clutch measures
270 from a single bird. A second example would be the incorrect use of residual variation
271 among pseudoreplicates to test the significance of fixed-effect variation at a different
272 level of a survey or experiment’s hierarchical design.

273 *Further reading: Harrison (2015) shows how poor replication of the random*
274 *intercept groups can give unstable model estimates. Zuur et al (2016) discuss the*
275 *importance of identifying dependency structures in your data.*

276 Deciding Model Structure for GLMMs

277 Choosing Error Structures and Link Functions

278 Linear models make various statistical assumptions, including additivity of the linear
279 predictors, independence of errors, equal variance of errors (homoscedasticity) and
280 Normality of errors (Gelman & Hill 2007 p. 46; Zuur et al 2009 p. 19). Ecologists often
281 deal with response variables that violate these assumptions, and face several decisions
282 about model specification to ensure models of such data are robust. The price for
283 ignoring violation of these assumptions tends to be an inflated Type I error rate (Zuur et
284 al 2010; Ives 2015).

285 For continuous response variables (e.g. mass, length), a Gaussian (also termed
286 Normal) error structure is often appropriate. Linear regression using a Gaussian
287 distribution will directly predict continuous data y from a linear predictor of covariates
288 (Lindsay 1974). Thus, model coefficients are on the same scale as the units of the
289 outcome variable e.g. mm of rainfall, or kg of mass. In some cases, however,
290 transformation of the response variable may still be required to improve the fit of a
291 Gaussian model. For example, the additivity assumption can be violated if there is a
292 non-linear relationship between the outcome variable and the predictors, but log-
293 transforming the outcome can often remedy this (Gelman & Hill 2007). Conversely, the
294 goal may be to quantify differences in mean mass between males and females, but if
295 the variance in mass for one sex is greater than the other, the assumption of
296 homogeneity of variance is violated. Transformation of the data can remedy this (Zuur
297 et al 2009), 'mean-variance stabilising transformations' ensure the variance around the
298 fitted mean of each group is similar, making the models more robust. Alternatively,
299 modern statistical tools such as the 'varIdent' function in the R package *nlme* can allow
300 one to explicitly model differences in variance between groups to avoid the need for
301 data transformation.

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

305

306 For non-Gaussian data, our modelling choices become more complex. Non-
307 Gaussian data structures include Poisson-distributed counts (number of eggs laid,
308 number of parasites); binomial-distributed constrained counts (number of eggs that
309 hatched in a clutch; prevalence of parasitic infection in a group of hosts) and Bernoulli-
310 distributed binary traits (e.g. infected with a parasite or not). Gaussian models of these
311 data would be inappropriate because they violate the assumptions of normality of errors
312 and homogenous variance. For example, the expected variance of a Poisson-distributed
313 variable is equal to its mean, so in a modelling context as the fitted mean increases so
314 too will the error around it. Binomial data has maximal variance at intermediate
315 probabilities, and zero variance when probabilities are zero or one. It is important to
316 mention, however, that real world data will only ever approximate a given distribution
317 and the correspondence of the data to the chosen distribution should be verified,
318 regardless of its 'type' (e.g count or proportion data). To model these data, we have two
319 initial choices: i) we can apply a transformation to our non-Gaussian response to 'make
320 it' approximately Gaussian, and then use a Gaussian model; or ii) we can apply a
321 GLMM and specify the appropriate error distribution and link function. The link function
322 takes into account the (assumed) empirical distribution of our data by transformation of
323 the linear predictor within the model, and so normalises the residuals of the model. It is
324 critical to note that transformation of the raw response variable is not equivalent to using
325 a link function to apply a transformation in the model. Data-transformation applies the
326 transformation to the raw response, whilst using a link function transforms the fitted
327 mean (the linear predictor). That is, the mean of a log-transformed response (using a
328 data transformation) is not identical to the logarithm of a fitted mean (using a link
329 function).

330 Crawley (2013 p. 560) gives the canonical link functions for the most common
331 generalized models: log for Poisson, logit (log-odds) for Binomial, and reciprocal for
332 Gamma errors. While it is beyond the scope of this paper to go through each possible
333 combination of error structure and link function, it is important to remember that several
334 combinations are possible depending on the structure of your data (see Zuur et al
335 2009). Your choice of link function may improve the fit of your model, but it is important
336 to know what assumptions your chosen link functions make about your data (Zuur et al

337 2009). The issue of transforming non-Gaussian data to fit Gaussian models to them is
338 contentious. Zuur et al (2009) suggest you should always use the appropriate modelling
339 tool e.g. Poisson GLMM for count data, or a generalized additive model (GAMM) for
340 non-linear data, rather than apply transformations just to be able to stay within the linear
341 modelling framework, as it can affect the influence of data points on the model (Keene
342 1995). For example, arcsin square-root transformation of proportion data was once
343 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
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.

355 An issue with transformations of non-Gaussian data is having to deal with zeroes
356 as special cases (e.g. you can't log transform a 0), so researchers often add a small
357 amount of noise to the zeroes to make the transformation work, a practice that has been
358 criticised (O'Hara & Kotze 2010). GLMMs remove the need for these 'adjustments' of
359 the data. The important point here is that transformations change the entire relationship
360 between Y and X (Zuur et al 2009), but different transformations do this to different
361 extents and it may be impossible to know which transformation is best without
362 performing simulations to test the efficacy of each (Warton & Hui 2011; Ives 2015).

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

368

369 Choosing Predictors and Interactions

370 One of the most important decisions during the modelling process is deciding which
371 predictors and interactions to include in models. Best practice demands that each model
372 should represent a specific *a priori* hypothesis concerning the drivers of patterns in data
373 (Burnham & Anderson 2002; Forstmeier & Schielzeth 2011), allowing you to assess the
374 relative support for these hypotheses in your data irrespective of model selection
375 philosophy. The definition of “hypothesis” must be broadened from the strict pairing of
376 null and alternative that is classically drilled into young pupils of statistics and
377 experimental design. Frequentist approaches to statistical modelling still work with
378 nested pairs of hypotheses. Information theorists work with whole sets of competing
379 hypotheses. Bayesian modellers are comfortable with the idea that every possible
380 parameter estimate is a hypothesis in its own right. But these epistemological
381 differences do not really help to solve the problem of “which” predictors should be
382 considered valid members of the full set to be used in a statistical modelling exercise. It
383 is therefore often unclear how best to design your most complex model, often referred
384 to as the *maximal model* (which contains all factors, interactions and covariates that
385 might be of any interest, Crawley 2013) or as the *global model* (a highly parameterized
386 model containing the variables and associated parameters thought to be important of
387 the problem at hand, Burnham & Anderson 2002; Grueber et al 2011). We shall use the
388 latter term here for consistency with terminology used in information-theory (Grueber et
389 al 2011).

390 Deciding which terms to include in the model requires careful and rigorous *a*
391 *priori* consideration of the system under study. This may appear obvious; however
392 diverse authors have noticed a lack of careful thinking when selecting variables for
393 inclusion in a model (Peters 1991, Chatfield 1995, Burnham & Anderson 2002). Lack of
394 *a priori* consideration, of what models represent, distinguishes rigorous hypothesis
395 testing from ‘fishing expeditions’ that seek significant predictors among a large group of
396 contenders. Ideally, the global model should be carefully constructed using the
397 researchers’ knowledge and understanding of the system such that only predictors likely

398 to be pertinent to the problem at hand are included, rather than including all the data the
399 researcher has collected and/or has available. This is a pertinent issue in the age of 'big
400 data', where researchers are often overwhelmed with predictors and risk skipping the
401 important step of *a priori* hypothesis design. In practice, for peer reviewers it is easy to
402 distinguish fishing expeditions from *a priori* hypothesis sets based on the evidence base
403 presented in introductory sections of research outputs.

404

405 *How Complex Should My Global Model Be?*

406 The complexity of the global model will likely be a trade-off between the number
407 of observations you have measured (the n of the study) and your proposed hypotheses
408 about how the measured variables affect the outcome (response) variable. Lack of
409 careful consideration of the parameters to be estimated can result in a model containing
410 more parameters than observations, called overparameterisation (Southwood &
411 Henderson 2000, Quinn & Keough 2002, Crawley 2013). In simple GLMs,
412 overparameterisation results in a rapid decline in (or absence of) degrees of freedom
413 with which to estimate residual error. Detection of overparameterisation in LMMs can be
414 more difficult because each random effect uses only a single degree of freedom,
415 however the estimation of variance among small numbers of groups can be numerically
416 unstable. Unfortunately, it is common practice to fit a global model that is simply as
417 complex as possible, irrespective of what that model actually represents; that is a
418 dataset containing k predictors yields a model containing a k -way interaction among all
419 predictors and simplify from there (Crawley 2013). This approach is flawed for two
420 reasons. First, this practice encourages fitting biologically-unfeasible models containing
421 nonsensical interactions. A good rule of thumb is that it should be possible to draw a
422 graph of what the fitted model 'looks like' for various combinations of predictors – failing
423 to draw the fitted lines of the 3-way interaction means refraining from fitting model
424 containing one. Second, using this approach makes it very easy to fit a model too
425 complex for the data. At best, the model will fail to converge, thus preventing inference.
426 At worst, the model will “work”, risking false inference. Guidelines for the ideal ratio of
427 data points (n) to estimated parameters (k) vary widely (see Forstmeier & Schielzeth
428 2011). Crawley (2013) suggests a minimum n/k of 3, though we argue this is very low

429 and that an n/k of 10 is more conservative. A 'simple' model containing a 3-way
430 interaction between continuous predictors and a single random intercept needs to
431 estimate 8 parameters, so requires a dataset of a *minimum* n of 80, and ideally >100 .
432 Interactions can be especially demanding, as fitting interactions between a multi-level
433 factor and a continuous predictor can result in poor sample sizes for specific treatment
434 combinations even if the total n is quite large (Zuur et al 2010), which will lead to
435 unreliable model estimates.

436 *Further reading: Zuur et al (2010) discuss data exploration techniques for*
437 *determining whether certain interactions should be included. Grueber et al (2011) show*
438 *an excellent worked example of a case where the most complex model is biologically*
439 *feasible and well-reasoned, containing only one 2-way interaction. Nakagawa and*
440 *Foster (2004) discuss the use of power analyses, which will be useful in determining the*
441 *appropriate n/k ratio for a given system.*

442

443 *Assessing Predictor Collinearity*

444 With the desired set of predictors identified, it is wise to check for collinearity among
445 predictor variables. Collinearity between predictors can cause several problems in
446 model interpretation because those predictors explain some of the same variance in
447 your response variable, and their effects cannot be estimated independently (Quinn and
448 Keough. 2002; Graham 2003): First, it can cause model convergence issues as models
449 struggle to partition variance between predictor variables. Second, positively correlated
450 variables can have negatively correlated regression coefficients, as the marginal effect
451 of one is estimated, given the effect of the other, leading to incorrect interpretations of
452 the direction of effects (Figure 2). Third, collinearity can inflate standard errors of
453 coefficient estimates and make 'true' effects harder to detect (Zuur et al 2010). Finally,
454 collinearity can affect the accuracy of model averaged parameter estimates during
455 multi-model inference (Freckleton 2011; Cade 2015). Examples of collinear variables
456 include climatic data such as temperature and rainfall, and morphometric data such as
457 body length and mass. Collinearity can be detected in several ways, including creating
458 correlation matrices between raw explanatory variables, with values >0.7 suggesting
459 both should not be used in the analysis (Dormann et al. 2013); or calculating the

460 variance inflation factor (VIF) of each predictor that is a candidate for inclusion in a
461 model (details in Zuur et al 2010) and dropping variables with a VIF higher than a
462 certain value (e.g. 3; Zuur et al 2010). One problem with these methods though is that
463 they rely on a user-selected, potentially arbitrary choice of threshold of either the
464 correlation coefficient or the VIF. Two solutions to this problem are to either select one
465 variable as representative of multiple collinear variables (Austin 2002), ideally using
466 biological knowledge/ reasoning to select the most meaningful variable (Zuur et al
467 2010); or conduct a dimension-reduction analysis (e.g. Principal Components Analysis;
468 James & McCullough 1990), leaving a single variable that accounts for most of the
469 shared variance among the correlated variables. Both approaches will only be
470 applicable if it is possible to group explanatory variables by common features, thereby
471 effectively creating broader, but still meaningful explanatory categories. For instance, by
472 using mass and body length metrics to create a 'scaled mass index' representative of
473 body size (Peig & Green 2009). In practice, any attempt to "tease apart" the relative
474 influence of two collinear predictors will fail. A common outcome is that the two
475 predictors of interest will share contribution to a principal component of the set of
476 predictors. This should be taken as strong indication that the predictors' signal cannot
477 be teased apart through inference, and experiments are required to manipulate them
478 independently.

479

480 *Standardising and Centering Predictors*

481 Transformations of predictor variables are common, and can improve model
482 performance and interpretability (Gelman & Hill 2007). Two common transformations
483 are i) predictor centering, where you subtract the mean of predictor x from every value
484 in x , giving a variable with mean 0 and SD on the original scale of x ; and ii) predictor
485 standardising, where you centre x but also divide by the SD of x , giving a variable with
486 mean 0 and SD 1. Rescaling the mean of predictors containing large values (e.g. rainfall
487 measured in thousands of mm) through centering/standardising will often solve
488 convergence problems, in part because the estimation of intercepts is brought into the
489 main body of the data themselves. Both approaches also remove the correlation
490 between main effects and their interactions, making main effects interpretable when

491 models also contain interactions (Schielzeth 2010). Note that this collinearity among
492 coefficients is distinct from collinearity between two separate predictors (see above).
493 Centering and standardising by the mean of a variable changes the interpretation of the
494 model intercept to the value of the outcome expected when x is at its mean value.
495 Standardising further adjusts the interpretation of the coefficient (slope) for x in the
496 model to the change in the outcome variable for a 1 SD change in the value of x .
497 Scaling is therefore a useful, indeed recommended, tool to improve the robustness of
498 regression models, but care must be taken in the interpretation and graphical
499 representation of outcomes.

500 *Further reading: Schielzeth (2010) provides an excellent reference to the*
501 *advantages of centering and standardising predictors. Gelman (2008) provides strong*
502 *arguments for standardising continuous variables by 2 SDs when you have binary*
503 *predictors in the model. Gelman & Hill (2007 p. 56, 434) discuss the utility of centering*
504 *by values other than the mean.*

505

506 **Quantifying GLMM Fit and Performance**

507 Once you have specified a global model, it is vital that you quantify model fit and report
508 these metrics in your manuscript to provide evidence that your model is robust. The
509 global model is considered the best candidate for assessing fit statistics such as
510 overdispersion (Burnham & Anderson 2002). Often, researchers will use information
511 criteria scores as a proxy for model fit, and claim that the large difference in AIC
512 between the top and null models is evidence of a good fit. This is incorrect: AIC tells us
513 nothing about whether the basic distributional and structural assumptions of the model
514 have been violated. Similarly a high R^2 value is in itself only a test of the magnitude of
515 model fit and not an adequate surrogate for proper model checks. Just because you
516 have a high R^2 value does not mean your model will pass checks for assumptions such
517 as homogeneity of variance. We strongly encourage researchers to view *model fit* and
518 *model adequacy* as two separate but equally important traits that must be assessed and
519 reported. Model fit can be poor for several reasons, including the presence of
520 overdispersion, failing to include interactions among predictors, failing to account for

521 non-linear relationships between variables, or specifying a sub-optimal error structure
522 and/or link function. Here we discuss some key metrics of fit and adequacy that should
523 be considered.

524

525 *Inspection of Residuals and Linear Model Assumptions*

526 Best practice is to examine plots of fitted values vs residuals for the entire model, as
527 well as model residuals versus all explanatory variables to look for patterns (Zuur et al
528 2010; 2016). In addition, there are further model checks specific to mixed models.
529 Firstly, you should inspect fitted values versus residuals for each grouping level of a
530 random intercept factor (Zuur et al 2009). This will often prove dissatisfying if there are
531 few data/residuals per group, however this in itself is a warning flag that the
532 assumptions of the model might be based on weak foundation. Another feature of fit
533 that is very rarely tested for in (G)LMMs is the assumption of normality of deviations of
534 the conditional modes of the random effects from the global intercept. Just as a
535 quantile-quantile (QQ) plot of linear model residuals should show points falling along a
536 straight line (e.g. Crawley 2007), so should a QQ plot of the random effect residuals.
537 Further reading: Zuur et al (2010) given an excellent overview of the assumptions of
538 linear models and how to test for their violation. See also Gelman & Hill (2007 p. 45).
539 The R package 'sjPlot' (Lüdtke 2017) has built in functions for several LMM
540 diagnostics, including random effect QQ plots. Zuur et al (2009) provides a vast
541 selection of model diagnostic techniques for a host of model types, including GLS,
542 GLMMs and GAMMS.

543

544 *Overdispersion*

545 If your model has a Gaussian (Normal) error structure, you should not be concerned
546 about overdispersion, as Gaussian models do not assume a specific mean-variance
547 relationship. For generalized mixed models (GLMMs) however (e.g. Poisson, Binomial),
548 the variance of the data can be greater than predicted by the error structure of your
549 model (e.g. Hilbe 2011). Overdispersion can be caused by several processes
550 influencing your data, including zero-inflation, aggregation (non-independence) among
551 counts, or both (Zuur et al 2009). The presence of overdispersion in your model

552 suggests it is a bad fit, and your parameter estimates and their standard errors will likely
553 be biased unless you account for the overdispersion (e.g. Harrison 2014). The use of
554 canonical binomial and Poisson error structures, when residuals are overdispersed,
555 tends to result in Type I errors because standard errors are underestimated. Adding an
556 observation-level random effect (OLRE) to overdispersed Poisson or Binomial models
557 can 'fix' the overdispersion and give more accurate estimates standard errors (Harrison
558 2014; 2015). However, OLRE models may yield inferior fit compared to models using
559 compound probability distributions such as the Negative-Binomial for count data (Hilbe
560 2011; Harrison 2014) or Beta-Binomial for proportion data (Harrison 2015), and so it is
561 good practice to assess the relative fit of both types of model using AIC before
562 proceeding (e.g. Zuur et al 2009). Researchers very rarely report the overdispersion
563 statistic (but see Elston et al 2001), but it should be made a matter of routine. See
564 'Assessing Model Fit Through Simulation' Section for advice on how to quantify and
565 model overdispersion.

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

571

572 R^2

573 In a linear modelling context, R^2 gives a measure of the proportion of explained variance
574 in the model, and is an intuitive metric for assessing model fit. Unfortunately, the issue
575 of calculating R^2 for (G)LMMs is particularly contentious; whereas for a simple linear
576 model with no random effects and a Normal error structure you can easily estimate the
577 residual variance, this is not the case for (G)LMMs. In fact, two issues exist with
578 generalising R^2 measures to (G)LMMs: i) for generalised models containing non-Normal
579 error structures, it is not clear how to calculate the residual variance term on which the
580 R^2 term is dependent; and ii) for mixed effects models, which are hierarchical in nature
581 and contain error (unexplained variance) at each of these levels, it is uncertain which
582 level to use to calculate a residual error term (Nakagawa & Schielzeth 2013). Diverse

583 methods have been proposed to account for this coefficient in GLMMs, including so-
584 called 'pseudo-r²' measures of explained variance (e.g. Nagelkerke 1991, Cox & Snell
585 1989), but their performance is often unstable for mixed models and can return negative
586 values (Nakagawa & Schielzeth 2013). Gelman & Pardoe (2006) derived a measure of
587 R² that accounts for the hierarchical nature of LMMs and gives a measure for both
588 group and unit level regressions (see also Gelman & Hill 2007 p. 474), but it was
589 developed for a Bayesian framework and a frequentist analogue does not appear to be
590 widely implemented. The method that has gained the most support over recent years is
591 that of Nakagawa & Schielzeth (2013).

592

593 The strength of the Nakagawa & Schielzeth (2013) method for GLMMs is that it returns
594 two complimentary R² values: the marginal R² encompassing variance explained by
595 only the fixed effects, and the conditional R² comprising variance explained by both
596 fixed and random effects i.e. the variance explained by the whole model (Nakagawa &
597 Schielzeth 2013). Ideally, both should be reported in publications as they provide
598 different information; which one is more 'useful' may depend on your rationale for
599 specifying random effects in the first instance. Note that when observation-level random
600 effects are included (see 'Overdispersion' section above), the conditional R² becomes
601 less useful as a measure of explained variance because it includes the extra-parametric
602 dispersion being modelled, but has no predictive power (Harrison 2014).

603 *Further reading: Nakagawa & Schielzeth (2013) provide an excellent and*
604 *accessible description of the problems with, and solutions to, generalising R² metrics to*
605 *GLMMs. The Nakagawa & Schielzeth (2013) R² functions have been incorporated into*
606 *several packages, including 'MuMIn' (Barton 2009) and 'piecewiseSEM' (Lefcheck*
607 *2015), and Johnson (2014) has developed an extension of the functions for random*
608 *slope models. See Harrison (2014) for a cautionary tale of how the GLMM R² functions*
609 *are artificially inflated for overdispersed models.*

610

611

612 *Stability of Variance Components and Testing Significance of Random Effects*

613 When models are too complex relative to the amount of data available, GLMM variance
614 components can collapse to zero (they cannot be negative). This is not a problem *per*
615 *se*, but it's important to acknowledge that in this case the model is equivalent to a
616 standard GLM. Reducing model complexity by removing interactions will often allow
617 random effects variance component estimates to become >0 , but this is problematic if
618 quantifying the interaction is the primary goal of the study. REML (restricted maximum
619 likelihood) should be used for estimating variance components of random effects in
620 Gaussian GLMMs as it produces less biased estimates compared to ML (maximum
621 likelihood) (Bolker et al 2009). However, when comparing two models with the same
622 random structure but different fixed effects, ML estimation cannot easily be avoided.
623 The RLRsim package (Scheipl, 2016) can be used to calculate restricted likelihood ratio
624 tests for variance components in mixed and additive models. Crucially, when testing the
625 significance of a variance component we are 'testing on the boundary' (Bolker et al
626 2009). That is the null hypothesis for random effects ($\sigma=0$) is at the boundary of its
627 possible range (it has to be ≥ 0), meaning p-values from a likelihood ratio test are
628 inaccurate. Dividing p values by 2 for tests of single variance components provides an
629 approximation to remedy this problem (Verbenke & Molenberghs, 2000).

630 Finally, estimating degrees of freedom for tests of random effects using Wald, t
631 or F tests or AICc is difficult, as a random effect can theoretically use anywhere
632 between 1 and $N - 1$ df (where N is the number of random-effect levels) (Bolker et al.
633 2009). Adequate F and P values can be calculated using Satterthwaite (1946)
634 approximations to determine denominator degrees of freedom implemented in the
635 package 'lmerTest' (Kuznetzova et al. 2014, see further details in section 'Model
636 Selection and Multi-Model Inference' below).

637

638 *Assessing Model Fit through Simulation*

639 Simulation is a powerful tool for assessing model fit (Gelman & Hill 2007; Kery 2010;
640 Zuur et al 2016), but is rarely used. The premise here is simple: for a given set of
641 parameter estimates (a model), if you were to simulate a dataset using those parameter
642 estimates, the fit of the model to those *simulated* 'ideal' data should be comparable to
643 the model's fit to the real data (Kery 2010). For each iteration, which yields a simulated

644 dataset, you can compute a statistic of interest such as the sum of squared residuals
645 (Kery 2010), the overdispersion statistic (Harrison 2014) or the percentage of zeroes for
646 a Poisson model (Zuur et al 2016). If the model is a good fit, after a sufficiently large
647 number of iterations (e.g. 10,000) the distribution of this test statistic should encompass
648 your observed statistic in the real data. Significant deviations outside of that distribution
649 indicate your model is a poor fit (Kery 2010). Figure 3 shows an example of using
650 simulation to assess the fit of a Poisson GLMM. After fitting a GLMM to count data, we
651 may wish to check for overdispersion and/or zero-inflation, the presence of which might
652 suggest we need to adjust our modelling strategy. Simulating 10,000 datasets from our
653 model reveals that the proportion of zeroes in our real data is comparable to simulated
654 expectation (Figure 3A). Conversely, simulating 1000 datasets and refitting our model to
655 each dataset, we see that the sum of the squared Pearson residuals for the real data is
656 far larger than simulated expectation (Figure 3B), giving evidence of overdispersion
657 (Harrison 2014). We can use the simulated frequency distribution of this test statistic to
658 derive a mean and 95% confidence interval for the overdispersion by calculating the
659 ratio of our test statistic to the simulated values (Harrison 2014). The dispersion statistic
660 for our model is 3.16 [95% CI 2.77 – 3.59]. Thus, simulations have allowed us to
661 conclude that our model is overdispersed, but that this overdispersion is not due to
662 zero-inflation. All R code for reproducing these simulations is provided in Online
663 Supplementary Material.

664 *Further reading: Rykiel (1996) discusses the need for validation of models in*
665 *ecology.*

666 Model Selection and Multi-Model Inference

667 Several methods of model selection are available once you have a robust global model
668 that satisfies standard assumptions of error structure and hierarchical independence
669 (Johnson & Omland 2004). We discuss the relative merits of each approach briefly
670 here, before expanding on the use of information-theory and multi-model inference in
671 ecology. We note that these discussions are not meant to be exhaustive comparisons,
672 and we encourage the reader to delve into the references provided for a comprehensive
673 picture of the arguments for and against each approach.

674

675 *Stepwise Selection, Likelihood Ratio Tests and P values*

676 A common approach to model selection is the comparison of a candidate model
677 containing a term of interest to the corresponding 'null' model lacking that term, using a
678 p value from a likelihood ratio test (LRT), referred to as null-hypothesis significance
679 testing (NHST; Nickerson 2000). Stepwise deletion involves using the NHST framework
680 to drop terms sequentially from the global model, and arrive at a 'minimal adequate
681 model' (MAM) containing only significant predictors (see Crawley 2013). NHST and
682 stepwise deletion have come under heavy criticism; they can overestimate the effect
683 size of 'significant' predictors (Whittingham et al 2006; Forstmeier & Schielzeth 2011)
684 and force the researcher to focus on a single best model as if it were the only
685 combination of predictors with support in the data. Although we strive for simplicity and
686 parsimony, this assumption is not reasonable in complex ecological systems (e.g.
687 Burnham, Anderson & Huyvaert 2011). It is common to present the MAM as if it arose
688 from a single *a priori* hypothesis, when in fact arriving at the MAM required multiple
689 significance tests (Whittingham et al 2006; Forstmeier & Schielzeth 2011). This cryptic
690 multiple testing can lead to hugely inflated Type I errors (Forstmeier & Schielzeth 2011).
691 Perhaps most importantly, LRT can be unreliable for fixed effects in GLMMs unless both
692 total sample size and replication of the random effect terms is high (see Bolker et al
693 2009 and references therein), conditions which are often not satisfied for most
694 ecological datasets. However, there are still cases where NHST may be the most
695 appropriate tool for inference. For example, in controlled experimental studies you may
696 wish to test the effect of a limited number of treatments and support estimates of effect
697 sizes with statements of statistical significance using model simplification (Mundry
698 2011). Importantly, Murtaugh (2009) found that the predictive ability of models assessed
699 using NHST was comparable to those selected using information-theoretic approaches
700 (see below), suggesting that NHST remains a valid tool for inference despite strong
701 criticism. Our advice is that NHST remains an important tool for analyses of
702 experiments and for inferential surveys with small numbers of well-justified *a priori*
703 hypotheses and with uncorrelated (or weakly correlated) predictors.

704 *Further reading: Stephens et al (2005) & Mundry (2011) argue the case for*
705 *NHST under certain circumstances such as well-designed experiments. Halsey et al*
706 *(2015) discuss the wider issues of the reliability of p values relative to sample size.*

707

708 *Information-Theory and Multi-Model Inference*

709 Unlike NHST, which leads you to focus on a single best model, model selection using
710 information theoretic (IT) approaches allows you to simultaneously rank the degree of
711 support in the data for several competing models using metrics such as Akaike's
712 Information Criterion (AIC). Information criteria attempt to quantify the Kullback-Leibler
713 distance (KLD), a measure of the relative amount of information lost when a given
714 model approximates the true data-generating process. Thus, relative difference among
715 models in AIC should be representative in relative differences in KLD, and the model
716 with the lowest AIC should lose the least information and be the best model in that it
717 optimises the trade-off between fit and complexity (e.g. Richards 2008). A key strength
718 of the IT approach is that it allows you to account for 'model selection uncertainty', the
719 idea that several competing models may all fit the data equally well (Burnham &
720 Anderson 2002; Burnham, Anderson & Huyvaert 2011). This is particularly useful when
721 competing models share equal "complexity" (i.e. number of predictors, or number of
722 residual degrees of freedom): in such situations, NHST is impossible because there is
723 no "null". Where several models have similar support in the data, inference can be
724 made from all models using model-averaging (Burnham & Anderson 2002; Johnson &
725 Omland 2004; Grueber et al 2011). Model averaging incorporates uncertainty by
726 weighting the parameter estimate of a model by that model's Akaike weight (often
727 referred to as the probability of that model being the best Kullback-Leibler model given
728 the data, but see Richards 2005). Multi-model inference places a strong emphasis on a
729 *priori* formulation of hypotheses (Burnham & Anderson 2002; Dochterman & Jenkins
730 2011; Lindberg et al 2015), and model-averaged parameter estimates arising from
731 multi-model inference are thought to lead to more robust conclusions about the
732 biological systems compared to NHST (Johnson & Omland 2004, but see Richards et al
733 2011). These strengths over NHST have meant that the use of IT approaches in
734 ecology and evolution has grown rapidly in recent years (Lindberg et al 2015; Barker &

735 Link 2015; Cade 2015). We do not expand on the specific details of the difference
736 between NHST and IT here, but point the reader to some excellent reference on the
737 topic. Instead, we use this section to highlight recent empirical developments in the best
738 practice methods for the application of IT in ecology and evolution.

739 *Further reading: Grueber et al (2011) and Symonds & Moussalli (2011) give a*
740 *broad overview of multi-model inference in ecology, and provide a worked model*
741 *selection exercise. Heygi & Garamszegi (2011) provide a detailed comparison of IT and*
742 *NHST approaches. Burnham, Anderson & Huyvaert (2011) demonstrate how AIC*
743 *approximates Kullback-Leibler information and provide some excellent guides for the*
744 *best practice of applying IT methods to biological datasets*

745

746 *Global Model Reporting*

747 Because stepwise deletion can cause biased effect sizes, presenting means and SEs of
748 parameters from the global model should be more robust, especially when the n/k ratio
749 is low (Forstmeier & Schielzeth 2011). A conservative approach relative to NHST is to
750 perform ‘full model tests’ (comparing the global model to an intercept only model) before
751 investigating single-predictor effects, as this controls the Type I error rate (Forstmeier &
752 Schielzeth 2011). Reporting the full model also helps reduce publication bias towards
753 strong effects, providing future meta-analyses with estimates of both significant and
754 non-significant effects (Forstmeier & Schielzeth 2011). Global model reporting should
755 not replace other model selection methods, but provides a robust measure of how likely
756 significant effects in minimal / best-AIC models are to arise by sampling variation alone.
757

758 **Practical Issues with Applying Information Theory to Biological Data**

759

760 *1. Using All-Subsets Selection or ‘Data Dredging’*

761 Dredging is the act of fitting a global model, often containing every possible interaction,
762 and then performing ‘all subsets’ selection on that model to fit every possible nested
763 model. On the surface, dredging might appear to be a convenient and fast way of
764 ‘uncovering’ the significant drivers of the patterns in your data. Dredging of enormous

765 global models containing large numbers of predictors and their interactions makes
766 analyses extremely prone to Type I errors and ‘overfitted’ models. Burnham & Anderson
767 (2002) caution strongly against dredging, and instead advocate ‘hard thinking’ about the
768 hypotheses underlying your data. If adopting an all subsets approach, it is worth noting
769 the number of models to consider increases exponentially with the number of predictors,
770 where 5 predictors require you to fit 2^5 (32) models, whilst 10 predictors requires 1024
771 models, both *without* including any interactions.

772 The inflation of Type I error rate through dredging is simple to demonstrate.
773 Figure 4 shows the results of a simulation exercise where we created datasets
774 containing various numbers of continuous and categorical variables, fitted a global
775 model containing all predictors as main effects and no interactions; and then dredged
776 that model. All simulated predictors were samples drawn from populations representing
777 the null hypothesis, i.e. having zero influence on the response variable. We considered
778 all models with an AIC score of within 6 of the best-supported AIC model to be equally
779 well supported (also referred to as the $\Delta 6$ AIC top model set, Richards 2008) (detailed
780 methods available in Online Supplementary Material). We assumed a Type I error had
781 occurred when the 95% confidence intervals for model averaged parameter estimates
782 from the $\Delta 6$ AIC set did not cross zero. The higher the number of terms in the model, the
783 higher the Type I error rate, reaching a maximum of over 60% probability of falsely
784 including a predictor in the top model set that was unrelated to the response variable.
785 Importantly, we found that the rate of increase (slope) in Type I error with added
786 continuous predictors was modified by the number of categorical variables (Fig. 4),
787 meaning the change in Type 1 error rate per continuous predictor was highest with
788 smaller numbers of categorical variables.

789 These results help to illustrate why dredging should not be common practice, and
790 you should not build global models containing huge numbers of variables and
791 interactions without prior thought about what the models represent for your study
792 system. In cases where you do perform all-subsets selection from a global model, it is
793 important to view these model selection exercises as exploratory (Symonds & Moussali
794 2011), and hold some data back from these exploratory analyses to be used for cross-
795 validation with your top model(s) (see Dochterman and Jenkins 2011 and references

796 therein). Here, you might use 90% of the data to fit the models and use the remaining
797 10% for confirmatory analysis to quantify how well the model(s) perform for prediction
798 (Zuur et al 2016). Such an approach requires a huge amount of data (Dochterman and
799 Jenkins 2011), but cross-validation to validate a model's predictive ability is rare and
800 should result in more robust inference (see also Fieberg & Johnson 2015).
801 Therefore, best practice is to consider only a handful of hypotheses and then build a
802 single statistical model to reflect each of them. This makes inference easier because the
803 resulting top model set will likely contain fewer parameters, and certainly fewer
804 spuriously 'significant' parameters (Burnham & Anderson 2002; Arnold 2010). However,
805 we argue 'dredging' may be sensible in a limited number of circumstances. For
806 example, if your most complex model contains two main effects and their interaction,
807 dredging that model will be indistinguishable from independently building the four
808 competing models nested in the global model, all of which may be considered likely to
809 be supported by the data. Similarly, when the global model is well-thought out, contains
810 few predictors, and only interactions likely to have empirical support, all-subsets
811 selection may be a valid variable selection tool. It is worth remembering that the Type I
812 error rate can quickly exceed the nominal 5% threshold if these conditions are not met
813 (Fig. 4). Moreover, a small number of models built to reflect well-reasoned hypotheses
814 are only valid if the predictors therein are not collinear (see 'Collinearity' section below).
815 All-subsets selection using the R package *MuMIn* (Barton 2016) will not automatically
816 check for collinearity, and so the onus falls on the researcher to be thorough in checking
817 for such problems.

818

819 2. *Deciding Which Information Criterion To Use*

820 Several information criteria are available to rank competing models, but their
821 calculations differ subtly. Commonly applied criteria include Akaike's Information
822 Criterion (AIC), the small sample size correction of AIC for when $n/k < 40$ (AICc), and the
823 Bayesian Information Criterion (BIC). QAIC is an adjustment to AIC that accounts for
824 overdispersion and should be used when you have identified overdispersion in your
825 model (see 'Overdispersion section' above). Note you do not have to use QAIC if you
826 have modelled the overdispersion in your dataset using zero-inflated models,

827 observation-level random effects, or compound probability distributions. Bolker et al
828 (2009) and Grueber et al (2011) provide details of how to calculate these criteria.

829 AIC maximises the fit/complexity trade-off of a model by balancing the model fit
830 with the number of estimated parameters. AICc and BIC both penalise the IC score
831 based on total sample size n , but the degree of penalty for AICc is less severe than BIC
832 for moderate sample sizes, and more severe for very low sample size (Brewer et al
833 2016). Whilst AIC tend to select overly complex models, Burnham and Anderson (2002)
834 criticised BIC for selecting overly simplistic models (underfitting). BIC is also criticised
835 because it operates on the assumption that the true model is in the model set under
836 consideration, whereas in ecological studies this is unlikely to be true (Burnham &
837 Anderson 2002; 2004). Issues exist with both AIC and BIC in a GLMM context for
838 estimating the number of parameters for a random effect (Bolker et al 2009; Grueber et
839 al 2011), and although df corrections to remedy this problem exist it is not always clear
840 what method is being employed by software packages (see Bolker et al 2009 Box 3).
841 Brewer et al (2016) show how the optimality of AIC, AICc and BIC for prediction
842 changes with both sample size and effect size of predictors (see also Burnham and
843 Anderson 2004). Therefore, the choice between the two metrics is not straightforward,
844 and may depend on the goal of the study i.e. model selection vs prediction, see Grueber
845 et al 2011 Box 1.

846

847 3. *Choice of Δ AIC Threshold*

848 Model averaging requires the identification of a “top model set” containing all models
849 with comparable support in the data, normally based on the change in AIC values
850 relative to the best AIC model (Δ AIC). Historically, Burnham & Anderson (2002)
851 recommended that only models with Δ AIC between 0-2 should be used for inference,
852 but subsequent work has shown that at least Δ 6 AIC is required to guarantee a 95%
853 probability that the best (expected) Kullback-Leibler Distance model is in the top model
854 set (Richards 2008; see also Burnham et al 2011). Alternatively, models can be ranked
855 by their Akaike weights and all those with an Akaike weight ≥ 0.95 retained in the “95%
856 confidence set” (Burnham & Anderson 2002), but doing so can lead to cumbersome top
857 model sets (Symonds & Moussali 2011). Using high cut-offs are not encouraged, to

858 avoid overly complex models followed by invalid results (Richards 2008; Grueber et al.
859 2011) but deciding on how many is too many remains a contentious issue (Grueber et
860 al. 2011). We suggest $\Delta 6$ as a minimum following Richards (2005; 2008).

861

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

863 It is well known that AIC tends towards overly complex models ('overfitting', Burnham &
864 Anderson 2002). As AIC only adds a 2 point penalty to a model for inclusion of a new
865 term, Arnold (2010) demonstrated that adding a nuisance predictor to a well-fitting
866 model leads to a ΔAIC value of the new model of ~ 2 , therefore appearing to warrant
867 inclusion in the top model set (see section above). Therefore, inference can be greatly
868 improved by eliminating models from the top model set that are more complex versions
869 of nested models with better AIC support, known as the nesting rule (Richards 2005;
870 2008; Richards, Whittingham & Stephens 2011). Doing so greatly reduces the number
871 of models to be used for prediction, and improves parameter accuracy (Arnold 2010;
872 Richards et al 2008). Symonds & Moussali (2011) caution that its applicability has not
873 yet been widely assessed over a range of circumstances, but the theory behind its
874 application is sound and intuitive (Arnold 2010). One potential problem is that once you
875 have removed models from the top model set, interpretation of the Akaike weights for
876 the remaining models becomes difficult, and thus model-averaged estimates using
877 these weights may not be sensible.

878

879 *5. Using Akaike Weights to Quantify Variable Importance*

880 Once you have arrived at a top model set, it is common practice to use the summed
881 Akaike weights of every model in that set in which a predictor of interest occurs as a
882 measure of 'variable importance' (e.g. Grueber et al 2011). Recent work has
883 demonstrated that this approach is flawed because Akaike weights are interpreted as
884 relative model probabilities, and give no information about the importance of individual
885 predictors in a model (Cade 2015). The sum of AIC weights as a measure of variable
886 importance may at best be a measure of how likely a variable would be included after
887 repeated sampling of the data (Burnham & Anderson 2002; Cade 2015). A better

888 measure of variable importance would be to compare standardised effect sizes
889 (Schielzeth 2010; Cade 2015).

890

891 *6. Model Averaging when Predictors Are Collinear*

892 The aim of model averaging is to incorporate the uncertainty of the size and presence of
893 effects among a set of candidate models with equal support in the data. Model
894 averaging using Akaike weights proceeds on the assumption that predictors are on
895 common scales across models and are therefore comparable. Unfortunately, the nature
896 of multiple regression means that the scale and sign of coefficients will change across
897 models depending on the presence or absence of other variables in a focal model
898 (Cade 2015). The issue of predictor scaling changing across models is particularly
899 exacerbated when predictors are collinear, even when VIF values are low (Burnham
900 and Anderson 2002; Lukacs, Burnham & Anderson 2010; Cade 2015). Cade (2015)
901 recommends standardising model parameters based on partial standard deviations to
902 ensure predictors are on common scales across models prior to model averaging
903 (details in Cade 2015). We stress again the need to assess multicollinearity among
904 predictors in multiple regression modelling before fitting models (Zuur et al 2016) and
905 before model-averaging coefficients from those models (Lukacs, Burnham & Anderson
906 2010; Cade 2015)

907

908

909 **Conclusion**

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

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

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

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

925 4. Data dredging or 'fishing expeditions' can be very risky and inflate the number of
926 false positives enormously. If you are going to include all combinations of predictors you
927 must have strong *a priori* justification.

928 5. If you still wish to include a large number of predictors, backwards selection and
929 NHST should be avoided, and ranking via AIC of all competing models is preferred. A
930 critical question that remains to be addressed is whether model selection based on
931 information theory is superior to NHST even in cases of balanced experimental designs
932 with few predictors.

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

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

939 8. Other modelling approaches such as Bayesian inference are available, and allow
940 much greater flexibility in choice of model structure, error structure and link function.
941 However, the ability to compare among competing models is underdeveloped, and
942 where these tools do exist (e.g. reversible-jump MCMC), they are not yet accessible
943 enough to non-experts to be useful.

944

945

946

947 **Acknowledgements**

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

950

951

952

953 References

- 954 Arnold TW. 2010. Uninformative parameters and model selection using Akaike's
955 Information Criterion. *Journal of Wildlife Management*, 74: 1175-1178.
- 956 Austin MP. 2002. Spatial prediction of species distribution: an interface between
957 ecological theory and statistical modelling. *Ecological Modelling* 157: 101–118.
- 958 Barker RJ, Link WA. 2015. Truth, models, model sets, AIC, and multimodel inference: A
959 Bayesian perspective. *The Journal of Wildlife Management*, 79: 730–738.
- 960 Bartoń K. 2016. MuMIn: Multi-Model Inference. R package version
961 1.15.6.<https://CRAN.R-project.org/package=MuMIn>
- 962 Bates D, Maechler M, Bolker B, Walker S. 2015. Fitting Linear Mixed-Effects Models
963 Using lme4. *Journal of Statistical Software* 67: 1-48.
- 964 Bolker BM, Brooks ME, Clark CJ, Geange SW, Poulsen JR, Stevens MHH, White JSS.
965 2009. Generalized linear mixed models: a practical guide for ecology and
966 evolution. *Trends in Ecology and Evolution* 24: 127–135.
- 967 Breslow NE, Clayton DG. 1993. Approximate inference in generalized linear mixed
968 models. *Journal of the American statistical Association* 88: 9-25.
- 969 Brewer MJ, Butler A, Cooksley SL. 2016. The relative performance of AIC, AICC and
970 BIC in the presence of unobserved heterogeneity. *Methods in Ecology and*
971 *Evolution*, 7: 679-692.
- 972 Burnham KP, Anderson DR. 2002. Model Selection and Multimodel Inference: A
973 Practical Information-Theoretic Approach, Second. Springer-Verlag, New York.
- 974 Burnham KP, Anderson DR. 2004. Multimodel inference: understanding AIC and BIC in
975 model selection. *Sociological Methods & Research*, 33: 261-304.
- 976 Burnham KP, Anderson DR, Huyvaert KP. 2011. AIC model selection and multimodel
977 inference in behavioral ecology: Some background, observations, and
978 comparisons. *Behavioral Ecology and Sociobiology*, 65: 23–35.
979 <http://doi.org/10.1007/s00265-010-1029-6>
- 980 Cade BS. 2015. Model averaging and muddled multimodel inferences. *Ecology* 96:
981 2370–2382.

- 982 Chatfield C. 1995. Model uncertainty, data mining and statistical inference (with
983 discussion). *Journal of the Royal Statistical Society, Series A* 158: 419-66.
- 984 Cox DR, Snell EJ. 1989. *The Analysis of Binary Data*, 2nd ed. London: Chapman and
985 Hall.
- 986 Crawley (2013) *The R Book*. Second Edition. Wiley, Chichester UK.
- 987 Dochtermann NA, Jenkins SH. 2011. Developing multiple hypotheses in behavioural
988 ecology. *Behavioral Ecology and Sociobiology* 65: 37-45.
- 989 Dominicus A, Skrondal A, Gjessing HK, Pedersen NL, Palmgren J. 2006. Likelihood ratio
990 tests in behavioral genetics: problems and solutions. *Behavior Genetics* 36: 331-
991 340
- 992 Dormann CF, Elith J, Bacher S, Buchmann C, Carl G, Carré G, Marquéz JR, Gruber B,
993 Lafourcade B, Leitão PJ, Münkemüller T. 2013. Collinearity: a review of methods
994 to deal with it and a simulation study evaluating their performance. *Ecography* 36:
995 027-046.
- 996 Ellison AM. 2004. Bayesian inference in ecology. *Ecology letters* 7: 509-520.
- 997 Elston, DA, Moss R, Boulinier T, Arrowsmith C, Lambin X, 2001. Analysis of
998 aggregation, a worked example: numbers of ticks on red grouse
999 chicks. *Parasitology* 122: 563-569.
- 1000 Fieberg J, Johnson DH. 2015. MMI: Multimodel inference or models with management
1001 implications? *The Journal of Wildlife Management* 79: 708-718.
- 1002 Forstmeier W, Schielzeth H. 2011. Cryptic multiple hypotheses testing in linear models:
1003 Overestimated effect sizes and the winner's curse. *Behavioral Ecology and*
1004 *Sociobiology*, 65: 47-55. <http://doi.org/10.1007/s00265-010-1038-5>
- 1005 Freckleton RP. 2011. Dealing with collinearity in behavioural and ecological data: model
1006 averaging and the problems of measurement error. *Behavioral Ecology and*
1007 *Sociobiology*, 65: 91-101.
- 1008 Garamszegi LZ. 2011. Information-theoretic approaches to statistical analysis in
1009 behavioural ecology: An introduction. *Behavioral Ecology and Sociobiology* 65: 1-
1010 11.
- 1011 Gelman A, Hill J. 2007. *Data analysis using regression and hierarchical/multilevel*
1012 *models*. New York, NY, USA: Cambridge University Press.

- 1013 Gelman A. 2008. Scaling regression inputs by dividing by two standard
1014 deviations. *Statistics in Medicine*, 27: 2865-2873.
- 1015 Gelman A, Pardoe I. 2006. Bayesian measures of explained variance and pooling in
1016 multilevel (hierarchical) models. *Technometrics*, 48: 241-251.
- 1017 Graham ME (2003) Confronting multicollinearity in multiple linear regression. *Ecology*
1018 84: 2809-2815
- 1019 Grueber CE, Nakagawa S, Laws RJ, Jamieson IG. 2011. Multimodel inference in
1020 ecology and evolution: Challenges and solutions. *Journal of Evolutionary Biology*
1021 24: 699–711.
- 1022 Harrison XA. 2014. Using observation-level random effects to model overdispersion in
1023 count data in ecology and evolution. *PeerJ* 2: e616.
- 1024 Harrison XA. 2015. A comparison of observation-level random effect and Beta-Binomial
1025 models for modelling overdispersion in Binomial data in ecology &
1026 evolution. *PeerJ*, 3: p.e1114.
- 1027 Halsey LG, Curran-Everett D, Vowler SL, Drummond GB. 2015. The fickle P value
1028 generates irreproducible results. *Nature Methods* 12: 179-185.
- 1029 Hegyi G, Garamszegi LZ. 2011. Using information theory as a substitute for stepwise
1030 regression in ecology and behaviour. *Behavioral Ecology and Sociobiology* 65: 69-
1031 76.
- 1032 Hilbe JM. 2011. *Negative binomial regression*. Cambridge University Press.
- 1033 Houslay T, Wilson A. 2017. Avoiding the misuse of BLUP in behavioral ecology.
1034 *Behavioral Ecology* arx023 doi:10.1093/beheco/arx023
- 1035 Ives AR. 2015. For testing the significance of regression coefficients, go ahead and
1036 log-transform count data. *Methods in Ecology and Evolution*, 6:, 828-835.
- 1037 James FC, McCullugh CF. 1990. Multivariate Analysis In Ecology And Systematics:
1038 Panacea Or Pandora Box. *Annual Review of Ecology and Systematics* 21: 129–
1039 166.
- 1040 Johnson JB, Omland KS. 2004. Model selection in ecology and evolution. *Trends in*
1041 *Ecology and Evolution*, 19: 101–108.
- 1042 Johnson PCD. 2014. Extension of Nakagawa & Schielzeth's R^2 GLMM to random
1043 slopes models. *Methods in Ecology and Evolution* 5: 944-946.

- 1044 Kass RE, Caffo BS, Davidian M, Meng XL, Yu B, Reid N. 2016. Ten simple rules for
1045 effective statistical practice. *PLoS computational biology*, 12: p.e1004961.
- 1046 Keene ON. 1995. The log transform is special. *Statistics in Medicine* 14: 811–819.
- 1047 Kéry M. 2010. Introduction to WinBUGS for ecologists: Bayesian approach to
1048 regression, ANOVA, mixed models and related analyses. Academic Press.
- 1049 Kuznetsova A, Brockhoff PB, Christensen RHB. 2014. Package ‘lmerTest’. Test for
1050 random and fixed effects for linear mixed effect models (lmer objects of lme4
1051 package). R package ver.2.
- 1052 Lefcheck JS. 2015. piecewiseSEM: Piecewise structural equation modeling in R for
1053 ecology, evolution, and systematics. *Methods in Ecology and Evolution*. 7: 573-
1054 579.<[doi:10.1111/2041-210X.12512](https://doi.org/10.1111/2041-210X.12512)>
- 1055 Lindberg MS, Schmidt JH, Walker J. 2015. History of multimodel inference via model
1056 selection in wildlife science. *The Journal of Wildlife Management* 79: 704–707.
- 1057 Lindsey JK. 1974. Construction and comparison of linear models. *Journal of the Royal*
1058 *Statistical Society. Series B* 36: 418 – 425.
- 1059 Low-Décarie E, Chivers C, Granados M. 2014. Rising complexity and falling explanatory
1060 power in ecology. *Frontiers in Ecology and the Environment* 12: 412-418.
- 1061 Lüdtke D. 2016. SjPlot: Data Visualization for Statistics in Social Science. 2016. *R*
1062 *package version*, 2(0).
- 1063 Lukacs PM, Burnham KP, Anderson DR. 2010. Model selection bias and Freedman’s
1064 paradox. *Annals of the Institute of Statistical Mathematics* 62: 117–125.
- 1065 Mundry R. 2011. Issues in information theory-based statistical inference—a
1066 commentary from a frequentist’s perspective. *Behavioral Ecology and*
1067 *Sociobiology*, 65: 57-68.
- 1068 Murtaugh PA. 2007. Simplicity and complexity in ecological data analysis. *Ecology* 88:
1069 56-62.
- 1070 Murtaugh PA. 2009. Performance of several variable-selection methods applied to real
1071 ecological data. *Ecology Letters* 10: 1061-1068.
- 1072 Nagelkerke NJ. 1991. A note on a general definition of the coefficient of determination.
1073 *Biometrika* 78: 691-692.

- 1074 Nakagawa S, Foster T. 2004. The case against retrospective statistical power analyses
1075 with an introduction to power analysis. *Acta Ethologica* 7: 103-108.
- 1076 Nakagawa S, Schielzeth H. 2010. Repeatability for Gaussian and non-Gaussian data: a
1077 practical guide for biologists. *Biological Reviews* 85: 935-956
- 1078 Nakagawa S, Schielzeth H. 2013. A general and simple method for obtaining R² from
1079 generalized linear mixed-effects models. *Methods in Ecology and Evolution* 4:
1080 133-142.
- 1081 Nickerson RS. 2000. Null Hypothesis Significance Testing: A Review of an Old and
1082 Continuing Controversy. *Psychological Methods* 5: 241-301.
- 1083 O'Hara RB, Kotze DJ. 2010. Do not log-transform count data. *Methods in Ecology and*
1084 *Evolution*, 1: 118-122.
- 1085 Peters RH. 1991. *A critique for ecology*. Cambridge University Press.
- 1086 Peig J, Green AJ. 2009. New perspectives for estimating body condition from
1087 mass/length data: the scaled mass index as an alternative method. *Oikos*, 118:
1088 1883-1891.
- 1089 Quinn GP, Keough MJ. 2002. *Experimental design and data analysis for biologists*.
1090 Cambridge University Press.
- 1091 R Core Team. 2016. R: A language and environment for statistical computing. R
1092 Foundation for Statistical Computing, Vienna, Austria. URL [https://www.R-](https://www.R-project.org/)
1093 [project.org/](https://www.R-project.org/).
- 1094 Richards SA. 2005. Testing ecological theory using the information-theoretic approach:
1095 examples and cautionary results. *Ecology*, 86: 2805-2814.
- 1096 Richards SA. 2008. Dealing with overdispersed count data in applied ecology. *Journal*
1097 *of Applied Ecology*, 45 218–227.
- 1098 Richards, SA, Whittingham MJ, Stephens PA. 2011. Model selection and model
1099 averaging in behavioural ecology: the utility of the IT-AIC framework. *Behavioral*
1100 *Ecology and Sociobiology*, 65: 77–89.
- 1101 Rykiel EJ. 1996. Testing ecological models: The meaning of validation. *Ecological*
1102 *Modelling* 90: 229-244.

- 1103 Scheipl F, Greven S, Kuechenhoff H. 2008. Size and power of tests for a zero random
1104 effect variance or polynomial regression in additive and linear mixed models.
1105 *Computational Statistics & Data Analysis* 52: 3283--3299.
- 1106 Schielzeth H, Forstmeier W. 2009. Conclusions beyond support: overconfident
1107 estimates in mixed models. *Behavioral Ecology* 20: 416-420
- 1108 Schielzeth H. 2010. Simple means to improve the interpretability of regression
1109 coefficients. *Methods in Ecology and Evolution* 1: 103-113
- 1110 Southwood TRE, Henderson PA. 2000. *Ecological methods*. John Wiley & Sons.
- 1111 Stephens PA, Buskirk SW, Hayward GD, Martinez Del Rio C. 2005. Information theory
1112 and hypothesis testing: a call for pluralism. *Journal of Applied Ecology* 42: 4-12.
- 1113 Symonds MRE, Moussalli A. 2011. A brief guide to model selection, multimodel
1114 inference and model averaging in behavioural ecology using Akaike's information
1115 criterion. *Behavioral Ecology and Sociobiology*, 65: 13-21.
- 1116 Vaida F, Blanchard S. 2005. Conditional Akaike information for mixed-effects models.
1117 *Biometrika* 92: 351-370
- 1118 Van de Pol M, Wright J. 2009. A simple method for distinguishing within-versus
1119 between-subject effects using mixed models. *Animal Behaviour* 77: 753-758.
- 1120 Verbenke G, Molenberghs G. 2000. Linear mixed models for longitudinal data. New
1121 York, Springer.
- 1122 Wilson AJ, Réale D, Clements MN, Morrissey MM, Postma E, Walling CA, Kruuk LEB,
1123 Nussey DH. 2010. An ecologist's guide to the animal model. *Journal of Animal*
1124 *Ecology*, 79: 13-26. doi:10.1111/j.1365-2656.2009.01639.x
- 1125 Warton D, Hui F. 2011. The arcsine is asinine: the analysis of proportions in ecology.
1126 *Ecology* 92: 3-10
- 1127 Wood SN, Goude Y, Shaw S. 2015. Generalized additive models for large data
1128 sets. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 64:,139-
1129 155.
- 1130 Whittingham MJ, Stephens PA, Bradbury RB, Freckleton RP 2006. Why do we still use
1131 stepwise modelling in ecology and behaviour? *Journal of Animal Ecology* 75:
1132 1182-1189.

- 1133 Zuur AF, Ieno EN, Walker NJ, Saveliev AA, Smith GM. 2009 *Mixed Effects Models and*
1134 *Extensions in Ecology with R* Springer, New York
- 1135 Zuur AF, Ieno EN, Elphick CS. 2010. A protocol for data exploration to avoid common
1136 statistical problems. *Methods in Ecology and Evolution*, 1: 3-14.
- 1137 Zuur AF, Ieno EN, 2016. A protocol for conducting and presenting results of
1138 regression-type analyses. *Methods in Ecology and Evolution*, 7: 636-645.
- 1139
- 1140

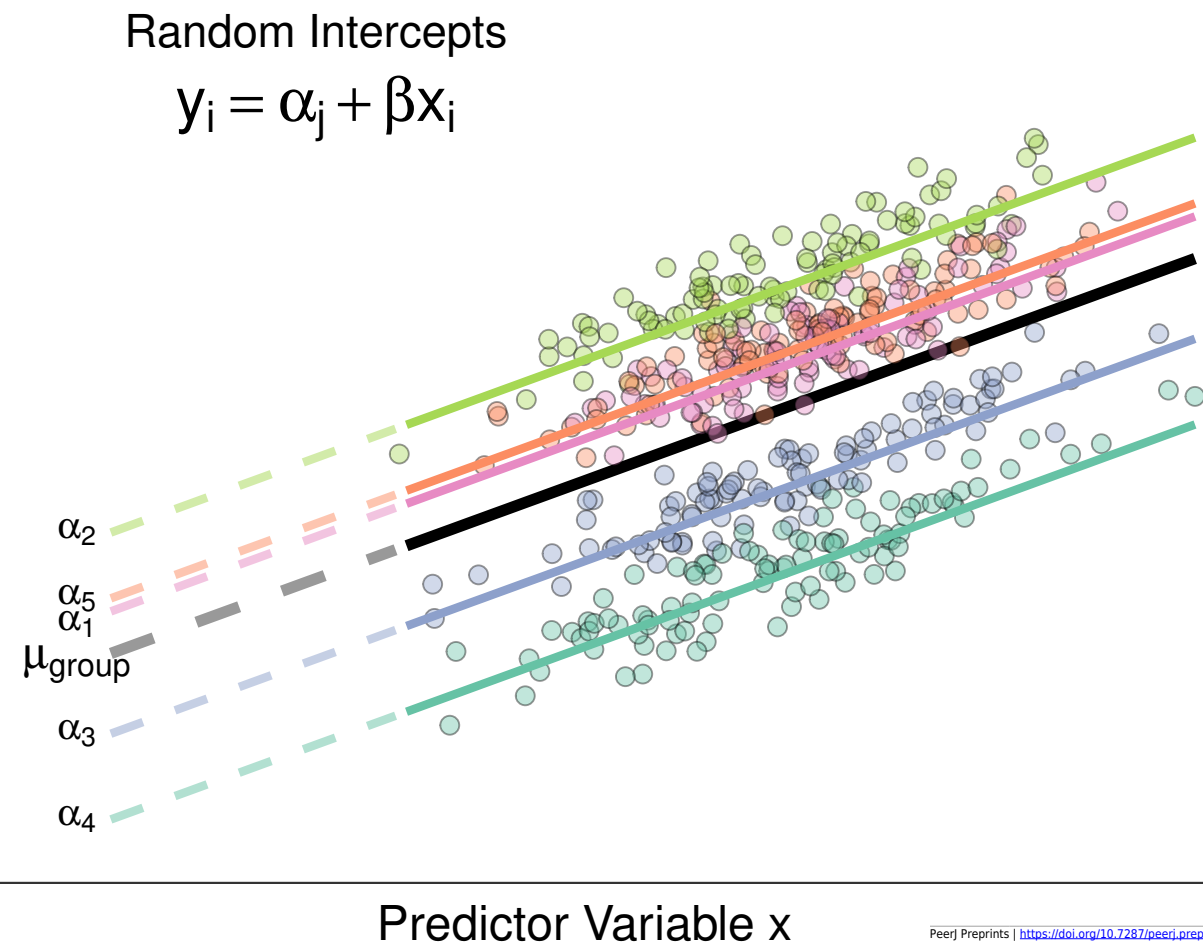
Figure 1(on next page)

Differences between Random Intercept vs Random Slope Models

(A) A random-intercepts model where the outcome variable y is a function of predictor x , with a random intercept for group ID (coloured lines). Because all groups have been constrained to have a common slope, their regression lines are parallel. Solid lines are the regression lines fitted to the data. Dashed lines trace the regression lines back to the y intercept (0 in this case). Point colour corresponds to group ID of the data point. The black line represents

A

Dependent Variable y

**B**

Dependent Variable y

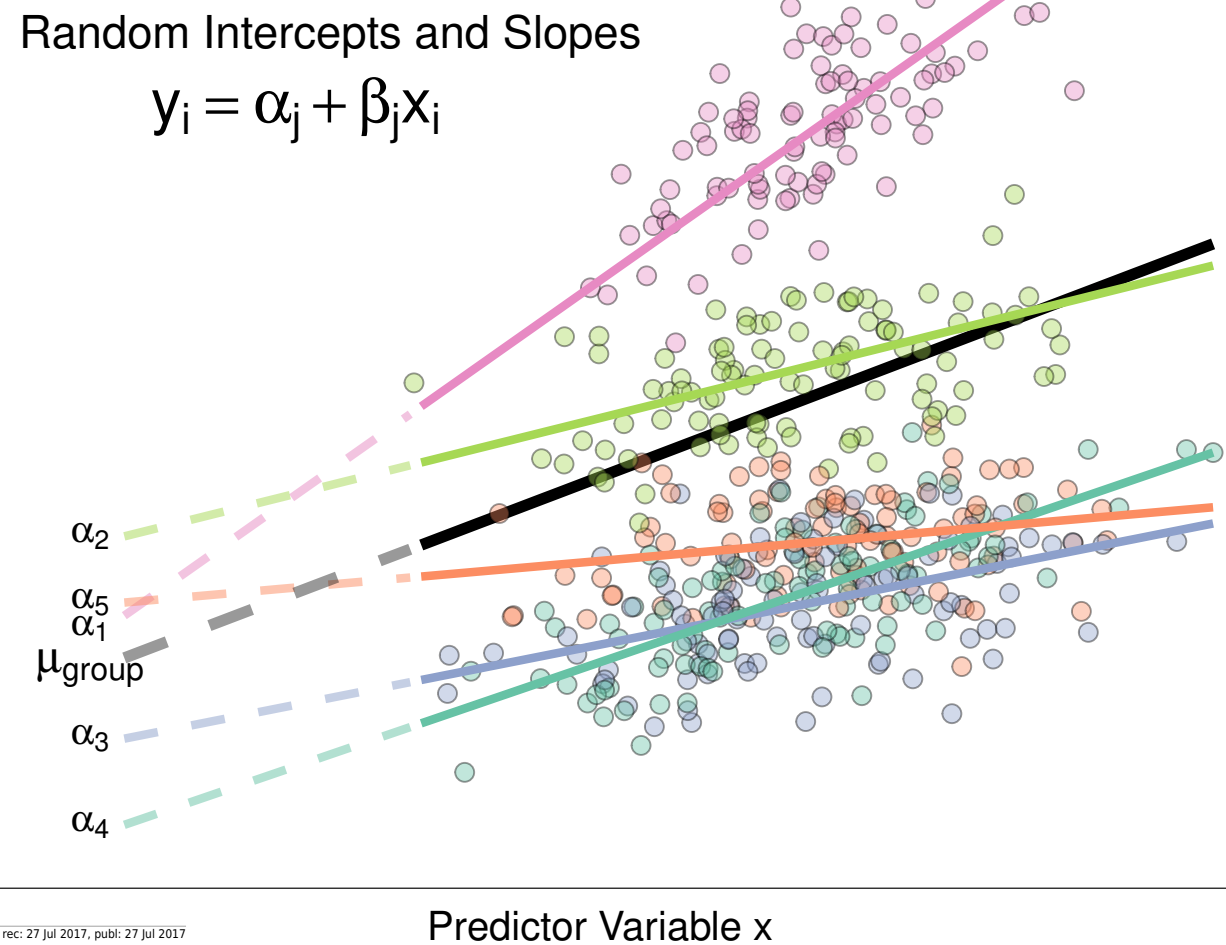


Figure 2(on next page)

The effect of collinearity on model parameter estimates.

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

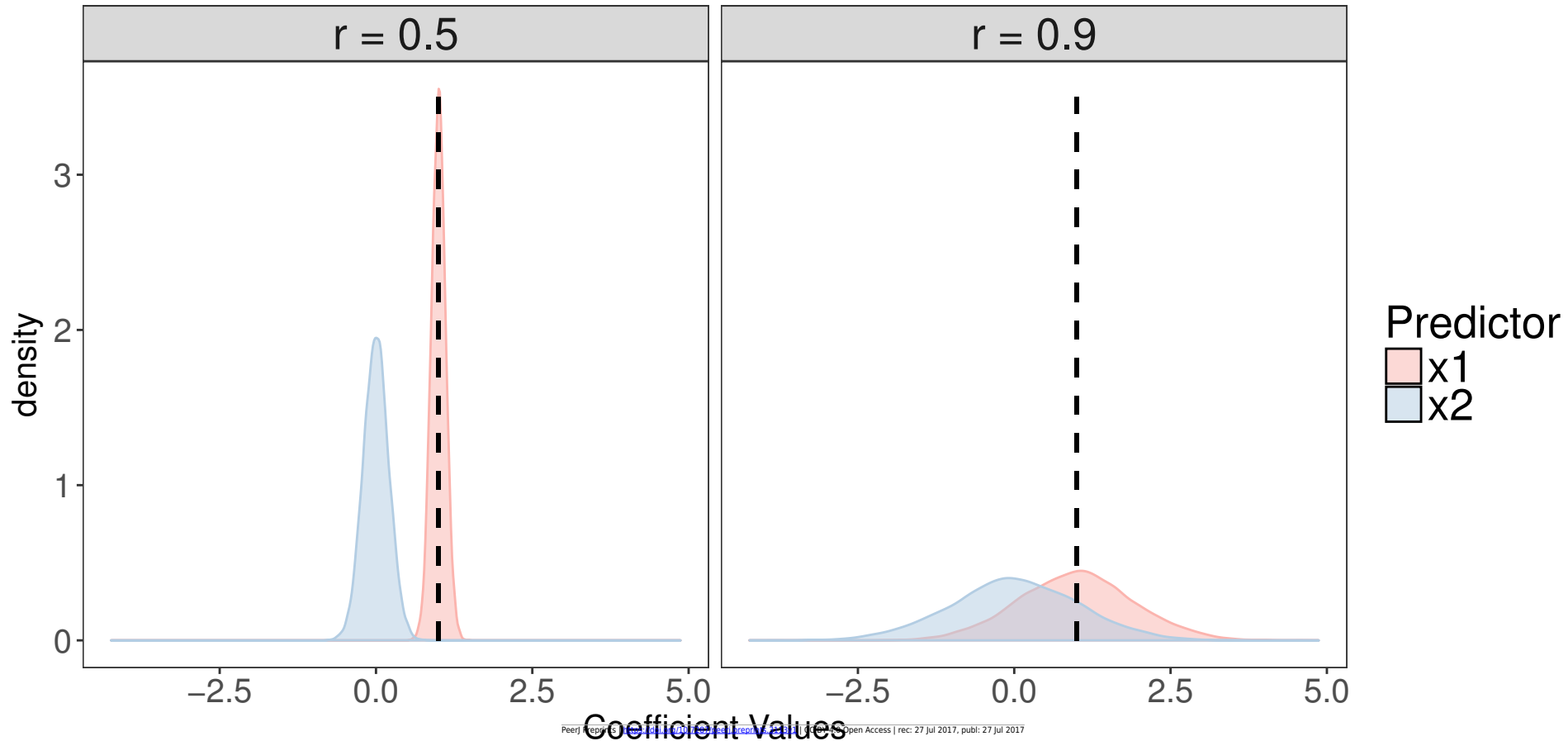


Figure 3(on next page)

Using Simulation to Assess Model Fit for GLMMs

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

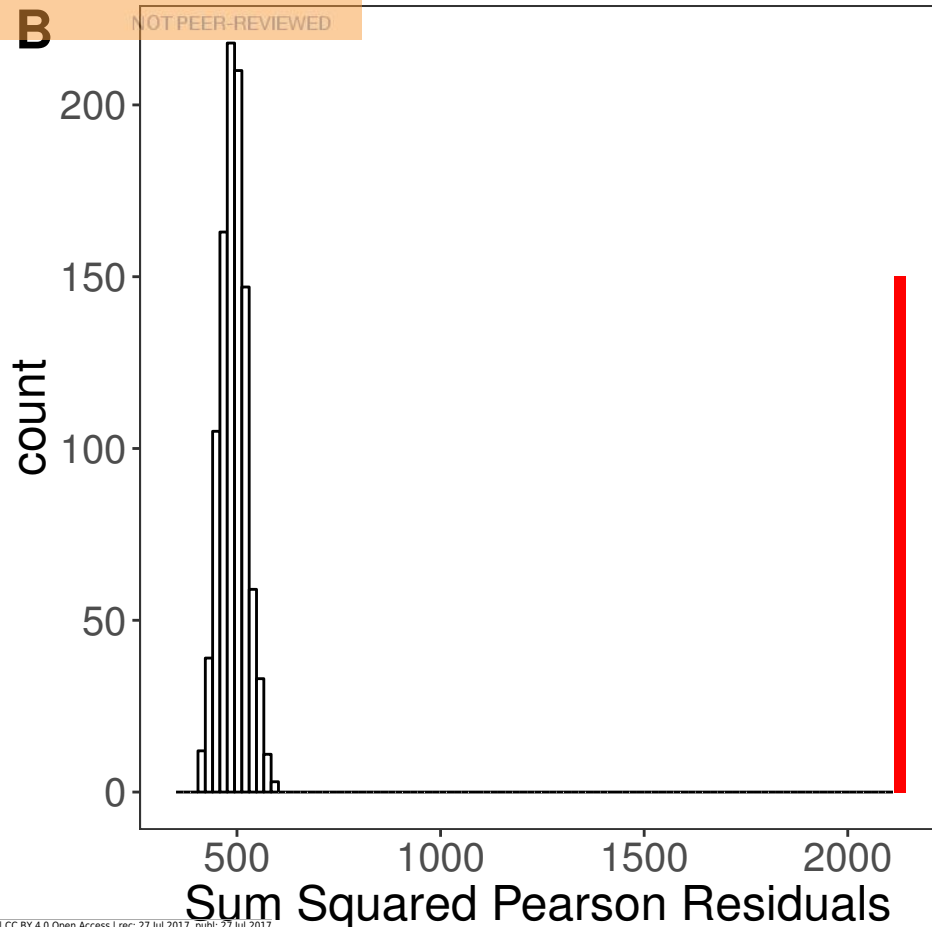
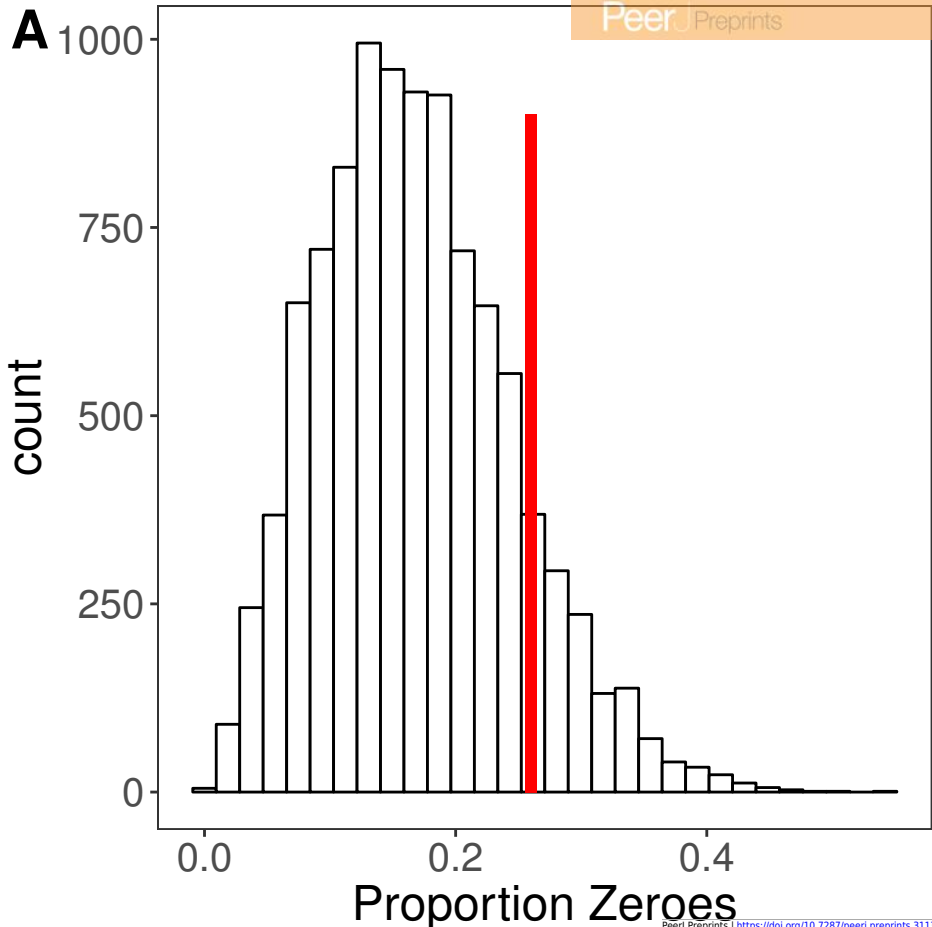


Figure 4(on next page)

The effect of data dredging on Type 1 Error Rate as a function of the number of continuous and categorical variables included in the global model

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

