# A brief introduction to mixed effects modelling and multi-model inference in ecology

**Xavier A Harrison** [Corresp., 1] , **Lynda Donaldson** [2, 3] , **Maria Eugenia Correa-Cano** [2] , **Julian Evans** [4, 5] , **David N Fisher** [4, 6] , **Cecily ED Goodwin** [2] , **Beth S Robinson** [2, 7] , **David J Hodgson** [4] , **Richard Inger** [2, 4]

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

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

[3] Wildfowl and Wetlands Trust, Slimbridge, Gloucestershire, United Kingdom

[4] Centre for Ecology and Conservation, University of Exeter, Penryn, United Kingdom

[5] Department of Biology, University of Ottawa, Ottawa, Canada

[6] Department of Integrative Biology, University of Guelph, Guelph, Canada

[7] WildTeam Conservation, Padstow, 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 regarding methods of model selection, with particular reference to the use of information theory and multi-model inference in ecology. 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    A Brief Introduction to Mixed Effects Modelling and Multi-model Inference in Ecology

2

3    Xavier A. Harrison[1], Lynda Donaldson[2&3], Maria Eugenia Correa-Cano[2], Julian

4    Evans[4,5,], David N. Fisher[4&6], Cecily E. D. Goodwin[2], Beth S. Robinson[2&7], David J.

5    Hodgson[4] and Richard Inger [2&4].

6

7    [1] Institute of Zoology, Zoological Society of London, London, United Kingdom

8    [2] Environment and Sustainability Institute, University of Exeter, Penryn, United Kingdom

9    [3]Wildfowl & Wetlands Trust, Slimbridge, Gloucestershire, United Kingdom

10    [4] Centre for Ecology and Conservation, University of Exeter, Penryn, United Kingdom

11    [5]Department of Biology, University of Ottawa, Ottawa, Canada

12    [6] Department of Integrative Biology, University of Guelph, Guelph, Canada

13    [7]WildTeam Conservation, Padstow, United Kingdom

14

15

16    Corresponding Authors:

17    Xavier Harrison xav.harrison@gmail.com

18    Richard Inger richinger@gmail.com

19

20

21

22

23

24

25

26

27

28

29

30

31

32    ABSTRACT

33

34    The use of linear mixed effects models (LMMs) is increasingly common in the analysis

35    of biological data. Whilst LMMs offer a flexible approach to modelling a broad range of

36    data types, ecological data are often complex and require complex model structures,

37    and the fitting and interpretation of such models is not always straightforward. The

38    ability to achieve robust biological inference requires that practitioners know how and

39    when to apply these tools. Here, we provide a general overview of current methods for

40    the application of LMMs to biological data, and highlight the typical pitfalls that can be

41    encountered in the statistical modelling process. We tackle several issues regarding

42    methods of model selection, with particular reference to the use of information theory

43    and multi-model inference in ecology. We offer practical solutions and direct the reader

44    to key references that provide further technical detail for those seeking a deeper

45    understanding. This overview should serve as a widely accessible guide to applying

46    LMMs to complex biological problems and model structures, and in doing so improve

47    the robustness of conclusions drawn from studies investigating ecological and

48    evolutionary questions.

49

50

51

## Introduction

53

54 In recent years, the suite of statistical tools available to biologists and the complexity of

55 biological data analyses have grown in tandem (Low-Decarie et al 2014; Zuur & Ieno

56 2016; Kass et al 2016). The availability of novel and sophisticated statistical techniques

57 means we are better equipped than ever to extract signal from noisy biological data, but

58 it remains challenging to know how to apply these tools, and which statistical

59 technique(s) might be best suited to answering specific questions (Kass et al 2016).

60 Often, simple analyses will be sufficient (Murtaugh 2007), but more complex data

61 structures often require more complex methods such as linear mixed effects models

62 (Zuur et al 2009), generalized additive models (Wood et al 2015) or Bayesian inference

63 (Ellison 2004). Both accurate parameter estimates and robust biological inference

64 require that ecologists be aware of the pitfalls and assumptions that accompany these

65 techniques and adjust modelling decisions accordingly (Bolker et al 2009).

66       Linear mixed effects models (LMMs) and generalized linear mixed effects models

67 (GLMMs), have increased in popularity in the last decade (Zuur et al 2009; Bolker et al

68 2009). Both extend traditional linear models to include a combination of fixed and

69 random effects as predictor variables. The introduction of random effects affords several

70 non-exclusive benefits. First, biological datasets are often highly structured, containing

71 clusters of non-independent observational units that are hierarchical in nature, and

72 LMMs allow us to explicitly model the non-independence in such data. For example, we

73 might measure several chicks from the same clutch, and several clutches from different

74 females, or we might take repeated measurements of the same chick's growth rate over

75 time. In both cases, we might expect that measurements within a statistical unit (here,

76 an individual, or a female's clutch) might be more similar than measurements from

77 different units. Explicit modelling of the random effects structure will aid correct

78 inference about fixed effects, depending on which level of the system's hierarchy is

79 being manipulated. In our example, if the fixed effect varies or is manipulated at the

80 level of the clutch, then treating multiple chicks from a single clutch as independent

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

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

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

112    with them. There have been several excellent papers in recent years on the use of

113    GLMMs in biology (Bolker et al 2009), the use of information theory and multi-model

114    inference for studies involving LMMs (Grueber et al 2011), best practice for data

115    exploration (Zuur et al 2009) and for conducting statistical analyses for complex

116    datasets (Zuur & Ieno 2016; Kass et al 2016). At the interface of these excellent guides

117    lies the theme of this paper: an updated guide for the uninitiated through the model

118    fitting and model selection processes when using LMMs. A secondary but no less

119    important aim of the paper is to bring together several key references on the topic of

120    LMMs, and in doing so act as a portal into the primary literature that derives, describes

121    and explains the complex modelling elements in more detail.

122        We provide a best practice guide covering the full analysis pipeline, from

123    formulating hypotheses, specifying model structure and interpreting the resulting

124    parameter estimates. The reader can digest the entire paper, or snack on each

125    standalone section when required. First, we discuss the advantages and disadvantages

126    of including both fixed and random effects in models. We then address issues of model

127    specification, and choice of error structure and/or data transformation, a topic that has

128    seen some debate in the literature (e.g. O'Hara & Kotze 2010; Ives 2015).  We also

129    address methods of model selection, and discuss the relative merits and potential

130    pitfalls of using information theory (IT), AIC and multi-model inference in ecology and

131    evolution. At all stages, we provide recommendations for the most sensible manner to

132    proceed in different scenarios. As with all heuristics, there may be situations where

133    these recommendations will not be optimal, perhaps because the required analysis or

134    data structure is particularly complex. If the researcher has concerns about the

135    appropriateness of a particular strategy for a given situation, we recommend that they

136    consult with a statistician who has experience in this area.


## 137   Understanding Fixed and Random Effects

138

139    A key decision of the modelling process is specifying model predictors as fixed or

140    random effects. Unfortunately, the distinction between the two is not always obvious,

141    and is not helped by the presence of multiple, often confusing definitions in the literature

142   (see Gelman and Hill 2007 p. 245). Absolute rules for how to classify something as a

143   fixed or random effect generally are not useful because that decision can change

144   depending on the goals of the analysis (Gelman and Hill 2007). We can illustrate the

145   difference between fitting something as a fixed (M1) or a random effect (M2) using a

146   simple example of a researcher who takes measurements of mass from 100 animals

147   from each of 5 different groups (n= 500) with a goal of understanding differences among

148   groups in mean mass. We use notation equivalent to fitting the proposed models in the

149   statistical software *R* (R Core Team 2016), with the LMMs fitted using the R package

150   *lme4* (Bates et al. 2015):

151

152                    M1 <- lm (mass ~ group)
153                   M2 <- lmer(mass ~ 1 + (1|group)

154

155   Fitting 'group' as a fixed effect in model M1 assumes the 5 'group' means are all

156   independent of one another, and share a common residual variance. Conversely, fitting

157   group as a random intercept model in model M2 assumes that the 5 measured group

158   means are only a subset of the realised possibilities drawn from a 'global' set of

159   population means that follow a Normal distribution with its own mean ($\mu_{group}$, Fig. 1A)

160   and variance ($\sigma^2_{group}$). Therefore, LMMs model the variance hierarchically, estimating

161   the processes that generate among-group variation in means, as well as variation within

162   groups. Treating groups from a field survey as only a subset of the *possible* groups that

163   could be sampled is quite intuitive, because there are likely many more groups (e.g.

164   populations) of the study species in nature than the 5 the researcher measured.

165   Conversely if one has designed an experiment to test the effect of three different

166   temperature regimes on growth rate of plants, specifying temperature treatment as a

167   fixed effect appears sensible because the experimenter has deliberately set the variable

168   at a given value of interest. That is, there are no unmeasured groups with respect to

169   that particular experimental design.

170         Estimating group means from a common distribution with known (estimated)

171   variance has some useful properties, which we discuss below, and elaborate on the

172  difference between fixed and random effects by using examples of the different ways

173  random effects are used in the literature.

174

175  *Controlling for non-independence among data points*

176  This is one of the most common uses of a random effect. Complex biological data sets

177  often contain nested and/or hierarchical structures such as repeat measurements from

178  individuals within and across units of time. Random effects allow for the control of non-

179  independence by constraining non-independent 'units' to have the same intercept

180  and/or slope (Zuur et al 2009; Zuur & Ieno 2016). Fitting *only* a random intercept allows

181  group means to vary, but assumes all groups have a common slope for a fitted

182  covariate (fixed effect). Fitting random intercepts *and* slopes allows the slope of a

183  predictor to vary based on a separate grouping variable. For example, one hypothesis

184  might be that the probability of successful breeding for an animal is a function of its

185  body mass. If we had measured animals from multiple sampling sites, we might wish to

186  fit 'sampling site' as a random intercept, and estimate a common slope (change in

187  breeding success) for body mass across all sampling sites by fitting it as a fixed effect:

188

189            M3 <- glmer(successful.breed ~ body.mass  +

190                (1|sample.site),family=binomial)

191

192  Conversely, we might wish to test the hypothesis that the strength of the effect (slope)

193  of body mass on breeding success varies depending on the sampling location i.e. the

194  change in breeding success for a 1 unit change in body mass is not consistent across

195  groups (Figure 1B). Here, 'body mass' is specified as a random slope by adding it to the

196  random effects structure. This model estimates a random intercept, random slope, and

197  the correlation between the two and also the fixed effect of body mass:

198

199            M4 <- glmer(successful.breed ~ body.mass +

200              (body.mass|sample.site),family=binomial)

201

202 Schielzeth & Forstmeier (2009); Barr et al (2013) and Aarts et al (2015) show that

203 constraining groups to share a common slope can inflate Type I and Type II errors.

204 Consequently, Grueber et al (2011) recommend always fitting both random slopes and

205 intercepts where possible. Whether this is feasible or not will depend on the data

206 structure (see 'Costs to Fitting Random Effects' section below). Figure 1 describes the

207 differences between random intercept models and those also containing random slopes.

208 *Further reading: Zuur & Ieno (2016) shows examples of the difficulties in*

209 *identifying the dependency structure of data and how to use flow charts / graphics to*

210 *help decide model structure. Kery (2010, Ch 12) has an excellent demonstration of how*

211 *to fit random slopes, and how model assumptions change depending on specification of*

212 *a correlation between random slopes and intercepts or not. Schielzeth & Forstmeier*

213 *(2009) and van de Pol & Wright (2009) are useful references for understanding the*

214 *utility of random slope models.*

215

216 *Improving the accuracy of parameter estimation*

217 Random effect models use data from all the groups to estimate the mean and variance

218 of the global distribution of group means. Assuming all group means are drawn from a

219 common distribution causes the estimates of their means to drift towards the global

220 mean $\mu_{group}$. This phenomenon, known as *shrinkage* (Gelman & Hill 2007; Kery 2010),

221 can also lead to smaller and more precise standard errors around means. Shrinkage is

222 strongest for groups with small sample sizes, as the paucity of within-group information

223 to estimate the mean is counteracted by the model using data from other groups to

224 improve the precision of the estimate. This 'partial pooling' of the estimates is a principal

225 benefit of fitting something as a random effect (Gelman & Hill 2007). However, it can

226 feel strange that group means should be shrunk towards the global mean, especially for

227 researchers more used to treating sample means as independent fixed effects.

228 Accordingly, one issue is that variance estimates can be hugely imprecise when there

229 are fewer than 5 levels of the random grouping variable (intercept or slope; see Harrison

230 2015). However, thanks to the Central Limit Theorem, the assumption of Gaussian

231 distribution of group means is usually a good one, and the benefits of hierarchical

232 analysis will outweigh the apparent costs of shrinkage.

233

234 *Estimating variance components*

235 In some cases, the variation among groups will be of interest to ecologists. For

236 example, imagine we had measured the clutch masses of 30 individual birds, each of

237 which had produced 5 clutches (n=150). We might be interested in asking whether

238 different females tend to produce consistently different clutch masses (high among-

239 female variance for clutch mass). To do so, we might fit an intercept-only model with

240 Clutch Mass as the response variable and a Gaussian error structure:

241

242
```
Model <- lmer(ClutchMass ~ 1 + (1|FemaleID)
```

243

244 By fitting individual 'FemaleID' as a random intercept term in the LMM, we estimate the

245 among-female variance in our trait of interest. This model will also estimate the residual

246 variance term, which we can use in conjunction with the among-female variance term to

247 calculate an 'intra-class correlation coefficient' that measures individual repeatability in

248 our trait (see Nakagawa & Schielzeth 2010). While differences among individuals can

249 be obtained by fitting individual ID as a fixed effect, this uses a degree of freedom for

250 each individual ID after the first, severely limiting model power, and does not benefit

251 from increased estimation accuracy through shrinkage. More importantly, repeatability

252 scores derived from variance components analysis can be compared across studies for

253 the same trait, and even across traits in the same study. Variance component analysis

254 is a powerful tool for partitioning variation in a focal trait among biologically interesting

255 groups, and several more complex examples exist (see Nakagawa & Schielzeth 2010;

256 Wilson et al 2010; Houslay & Wilson 2017). In particular, quantitative genetic studies

257 rely on variance component analysis for estimating the heritability of traits such as body

258 mass or size of secondary sexual characteristics (Wilson et al 2010). We recommend

259 the tutorials in Wilson et al (2010) and Houslay & Wilson (2017) for a deeper

260 understanding of the power and flexibility of variance component analysis.

261

262 *Making predictions for unmeasured groups*

263   Fixed effect estimates prevent us from making predictions for new groups because the
264   model estimates are only relevant to groups in our dataset (e.g. Zuur et al 2009 p. 327).
265   Conversely, we can use the estimate of the global distribution of population means to
266   predict for the average group using the mean of the distribution $\mu_{group}$ for a random
267   effects model (see Fig. 1). We could also sample hypothetical groups from our random
268   effect distribution, as we know its mean and SD (Zuur & Ieno 2016). Therefore, whether
269   something is fitted as a fixed or random effect can depend on the goal of the analysis:
270   are we only interested in the mean values for each group in our dataset, or do we wish
271   to use our results to extend our predictions to new groups? Even if we do not want to
272   predict to new groups, we might wish to fit something as a random effect to take
273   advantage of the shrinkage effect and improved parameter estimation accuracy.
274
275   **Considerations When Fitting Random Effects**
276   Random effect models have several desirable properties (see above), but their use
277   comes with some caveats. First, they are quite 'data hungry'; requiring at least 5 'levels'
278   (groups) for a random intercept term to achieve robust estimates of variance (Gelman &
279   Hill 2007; Harrison 2015). With <5 levels, the mixed model may not be able to estimate
280   the among-population variance accurately. In this case, the variance estimate will either
281   collapse to zero, making the model equivalent to an ordinary GLM (Gelman & Hill 2007
282   p. 275) or be non-zero but incorrect if the small number of groups that were sampled
283   are not representative of true distribution of means (Harrison 2015). Second, models
284   can be unstable if sample sizes across groups are highly unbalanced i.e. if some groups
285   contain very few data. These issues are especially relevant to random slope models
286   (Grueber et al 2011). Third, an important issue is the difficulty in deciding the
287   "significance" or "importance" of variance among groups. The variance of a random
288   effect is inevitably at least zero, but how big does it need to be to be considered of
289   interest? Fitting a factor as a fixed effect provides a statement of the significance of
290   differences (variation) among groups relatively easily. Testing differences among levels
291   of a random effect is made much more difficult for frequentist analyses, though not so in
292   a Bayesian framework (Kery 2010, see '*Testing Significance of Random Effects*'
293   section). Finally, an issue that is not often addressed is that of mis-specification of

294   random effects. GLMMs are powerful tools, but incorrectly parameterising the random

295   effects in the model could yield model estimates that are as unreliable as ignoring the

296   need for random effects altogether. Examples include: i) failure to recognise non-

297   independence caused by nested structures in the data e.g. multiple clutch measures

298   from a single bird; ii) failing to specify random slopes to prevent constraining slopes of

299   predictors to be identical across clusters in the data (see Barr et al 2013); and iii) testing

300   the significance of fixed effects at the wrong 'level' of hierarchical models that ultimately

301   leads to pseudoreplication and inflated Type I error rates. Traditionally users of LMMs

302   might have used *F*-tests of significance. *F*-tests are ill-advised for unbalanced

303   experimental designs and irrelevant for non-Gaussian error structures, but they at least

304   provide a check of model hierarchy using residual degrees of freedom for fixed effects.

305   The now-standard use of likelihood ratio tests of significance in LMMs means that users

306   and readers have little opportunity to check the position of significance tests in the

307   hierarchy of likelihoods.

308        *Further reading: Harrison (2015) shows how poor replication of the random*

309   *intercept groups can give unstable model estimates. Zuur & Ieno (2016) discuss the*

310   *importance of identifying dependency structures in the data.*


# Deciding Model Structure for GLMMs

311


### Choosing Error Structures and Link Functions

312

313   General linear models make various statistical assumptions, including additivity of the

314   linear predictors, independence of errors, equal variance of errors (homoscedasticity)

315   and normality of errors (Gelman & Hill 2007 p. 46; Zuur et al 2009 p. 19). Ecologists

316   often deal with response variables that violate these assumptions, and face several

317   decisions about model specification to ensure models of such data are robust. The price

318   for ignoring violation of these assumptions tends to be an inflated Type I error rate (Zuur

319   et al 2010; Ives 2015). In some cases, however, transformation of the response variable

320   may be required to ensure these assumptions are met. For example, an analytical goal

321   may be to quantify differences in mean mass between males and females, but if the

322   variance in mass for one sex is greater than the other, the assumption of homogeneity

323 of variance is violated. Transformation of the data can remedy this (Zuur et al 2009);

324 'mean-variance stabilising transformations' aim to make the variance around the fitted

325 mean of each group homogenous, making the models more robust. Alternatively,

326 modern statistical tools such as the 'varIdent' function in the R package *nlme* can allow

327 one to explicitly model differences in variance between groups to avoid the need for

328 data transformation.

329 *Further reading: Zuur et al (2010) provide a comprehensive guide on using data*

330 *exploration techniques to check model assumptions, and give advice on*

331 *transformations.*

332

333       For non-Gaussian data, our modelling choices become more complex. Non-

334 Gaussian data structures include, for example, Poisson-distributed counts (number of

335 eggs laid, number of parasites); binomial-distributed constrained counts (number of

336 eggs that hatched in a clutch; prevalence of parasitic infection in a group of hosts) or

337 Bernoulli-distributed binary traits (e.g. infected with a parasite or not). Gaussian models

338 of these data would violate the assumptions of normality of errors and homogenous

339 variance. To model these data, we have two initial choices: i) we can apply a

340 transformation to our non-Gaussian response to 'make it' approximately Gaussian, and

341 then use a Gaussian model; or ii) we can apply a GL(M)M and specify the appropriate

342 error distribution and link function. The link function takes into account the (assumed)

343 empirical distribution of our data by transformation of the linear predictor within the

344 model. It is critical to note that transformation of the raw response variable is not

345 equivalent to using a link function to apply a transformation in the model. Data-

346 transformation applies the transformation to the raw response, whilst using a link

347 function transforms the fitted mean (the linear predictor). That is, *the mean of a log-*

348 *transformed response (using a data transformation) is not identical to the logarithm of a*

349 *fitted mean (using a link function)*.

350       The issue of transforming non-Gaussian data to fit Gaussian models to them is

351 contentious. For example, arcsin square-root transformation of proportion data was

352 once extremely common, but recent work has shown it to be unreliable at detecting real

353 effects (Warton & Hui 2011). Both logit-transformation (for proportional data) and

354 Binomial GLMMs (for binary response variables) have been shown to be more robust

355 (Warton & Hui 2011). O'Hara & Kotze (2010) argued that log-transformation of count

356 data performed well in only a small number of circumstances (low dispersion, high

357 mean counts), which are unlikely to be applicable to ecological datasets. However, Ives

358 (2015) recently countered these assumptions with evidence that transformed count data

359 analysed using LMMs can often outperform Poisson GLMMs. We do not make a case

360 for either here, but acknowledge the fact that there is unlikely to be a universally best

361 approach; each method will have its own strengths and weakness depending on the

362 properties of the data (O'Hara & Kotze 2010). Checking the assumptions of the LMM or

363 GLMM is an essential step (see section 'Quantifying GLMM Fit and Performance'). An

364 issue with transformations of non-Gaussian data is having to deal with zeroes as special

365 cases (e.g. you can't log transform a 0), so researchers often add a small constant to all

366 data to make the transformation work, a practice that has been criticised (O'Hara &

367 Kotze 2010). GLMMs remove the need for these 'adjustments' of the data. The

368 important point here is that transformations change the entire relationship between Y

369 and X (Zuur et al 2009), but different transformations do this to different extents and it

370 may be impossible to know which transformation is best without performing simulations

371 to test the efficacy of each (Warton & Hui 2011; Ives 2015).

372 *Further reading: Crawley (2013 Ch 13) gives a broad introduction to the various error*

373 *structures and link functions available in the R statistical framework. O'Hara & Kotze*

374 *(2010); Ives (2015) and Warton et al (2016) argue the relative merits of GLMs vs log-*

375 *transformation of count data; Warton & Hui (2011) address the utility of logit-*

376 *transformation of proportion data compared to arcsin square-root transformation.*

377

378 **Choosing Random Effects I: Crossed or Nested?**

379 A common issue that causes confusion is this issue of specifying random effects as

380 either 'crossed' or 'nested'. In reality, the way you specify your random effects will be

381 determined by your experimental or sampling design (Schielzeth & Nakagawa 2013). A

382 simple example can illustrate the difference. Imagine a researcher was interested in

383 understanding the factors affecting the clutch mass of a passerine bird. They have a

384    study population spread across 5 separate woodlands, each containing 30 nest boxes.

385    Every week during breeding they measure the foraging rate of females at feeders, and

386    measure their subsequent clutch mass. Some females have multiple clutches in a

387    season and contribute multiple data points. Here, female ID is said to be *nested within*

388    *woodland*: each woodland contains multiple females unique to that woodland (that

389    never move among woodlands). The nested random effect controls for the fact that i)

390    clutches from the same female are not independent, and ii) females from the same

391    woodland may have clutch masses more similar to one another than to females from

392    other woodlands

393

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

395

396    Now imagine that this is a long-term study, and the researcher returns every year for 5

397    years to continue with measurements. Here it is appropriate fit year as a *crossed*

398    random effect because every woodland appears multiple times in every year of the

399    dataset, and females that survive from one year to the next will also appear in multiple

400    years.

401

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

403

404    Understanding whether your experimental/sampling design calls for nested or crossed

405    random effects is not always straightforward, but it can help to visualise experimental

406    design by drawing it (see Schielzeth and Nakagawa 2013 Fig. 1), or tabulating your

407    observations by these grouping factors (e.g. with the '*table*' command in R) to identify

408    how your data are distributed. We advocate that researchers always ensure that their

409    levels of random effect grouping variables are uniquely labelled. For example, females

410    are labelled 1 – *n* in each woodland, the model will try and pool variance for all females

411    with the same code. Giving all females a unique code makes the nested structure of the

412    data is implicit, and a model specified as ~ (1| Woodland) + (1|FemaleID) would be

413    identical to the model above.

414       Finally, we caution that whether two factors are nested or crossed affects the

415    ability of (G)LMMs to estimate the effect of the interaction between those two factors on

416    the outcome variable. Crossed factors allow the model to accurately estimate the

417    interaction effects between the two, whereas nested factors automatically pool those

418    effects in the second (nested) factor (Schielzeth and Nakagawa 2013). We do not

419    expand on this important issue here but direct the reader to Schielzeth and Nakagawa

420    2013 for an excellent treatment of the topic.

421    **Choosing Random Effects II: Random Slopes**

422    Fitting random slope models in ecology is not very common. Often, researchers fit

423    random intercepts to control for non-independence among measurements of a statistical

424    group (e.g. birds within a woodland), but force variables to have a common slope across

425    all experimental units. However, there is growing evidence that researchers should be

426    fitting random slopes as standard practice in (G)LMERs. Random slope models allow

427    the coefficient of a predictor to vary based on clustering / non-independence in the data

428    (see Fig. 1B). In our bird example above, we might fit a random slope for the effect of

429    foraging rate on clutch mass given each individual bird ID. That is, the magnitude of the

430    effect foraging rate on resultant clutch mass differs among birds. Random slope models

431    (also often called random coefficients models, Kery 2010) apply to both continuous and

432    factor variables. For example, if we had applied a two-level feeding treatment to birds in

433    each woodland (vitamin supplementation or control), we might also expect the

434    magnitude of the effect of receiving vitamin supplementation to differ depending on

435    which woodland it was applied to. So here we would specify random slopes for the

436    treatment variable given woodland ID.

437

438    Schielzeth & Forstmeier (2009) found that including random slopes controls Type I error

439    rate (yields more accurate p values), but also gives more power to detect among

440    individual variation. Barr et al (2013) suggest that researchers should fit the maximal

441    random effects structure possible for the data. That is, if there are four predictors under

442    consideration, all four should be allowed to have random slopes. However, we believe

443    this is unrealistic because random slope models require large numbers of data to

444  estimate variances and covariances accurately (Bates et al 2015). Ecological datasets

445  can often struggle to estimate a single random slope, diagnosed by a perfect correlation

446  (1 or -1) between random intercepts and slopes (Bates et al 2015). Therefore, the

447  approach of fitting the 'maximal' complexity of random effects structure (Barr et al 2013)

448  is perhaps better phrased as fitting the most complex mixed effects structure allowed by

449  your data (Bates et al 2015), which may mean either i) fitting random slopes but

450  removing the correlation between intercepts and slopes; or ii) fitting no random slopes

451  at all but accepting that this likely inflates the Type I error rate (Schielzeth & Forstmeier

452  2009). If fitting a random slope model including correlations between intercepts and

453  slopes, always inspect the intercept-slope correlation coefficient in the

454  variance/covariance summary returned by packages like *lme4* to look for evidence of

455  perfect correlations, indicative of insufficient data to estimate the model.

456  *Further Reading: Forstmeier and Schielzeth (2009) is essential reading for*

457  *understanding how random slopes control Type I error rate, and Bates et al (2015)*

458  *gives sound advice on how to iteratively determine optimal complexity of random effect*

459  *structure. Martin et al. (2011) conducted a simulation study identifying the sample sizes*

460  *necessary to accurately estimate parameters in random slope models. Barr et al (2013)*

461  *and Aarts et al (2015) discuss the merits of fitting random slopes to clustered data to*

462  *control false positive rates.*

463  **Choosing Fixed Effect Predictors and Interactions**

464  One of the most important decisions during the modelling process is deciding which

465  predictors and interactions to include in models. Best practice demands that each model

466  should represent a specific *a priori* hypothesis concerning the drivers of patterns in data

467  (Burnham & Anderson 2002; Forstmeier & Schielzeth 2011), allowing the assessment of

468  the relative support for these hypotheses in the data irrespective of model selection

469  philosophy. The definition of "hypothesis" must be broadened from the strict pairing of

470  null and alternative that is classically drilled into young pupils of statistics and

471  experimental design. Frequentist approaches to statistical modelling still work with

472  nested pairs of hypotheses. Information theorists work with whole sets of competing

473  hypotheses. Bayesian modellers are comfortable with the idea that every possible

474    parameter estimate is a hypothesis in its own right. But these epistemological

475    differences do not really help to solve the problem of "which" predictors should be

476    considered valid members of the full set to be used in a statistical modelling exercise. It

477    is therefore often unclear how best to design the most complex model, often referred to

478    as the *maximal model* (which contains all factors, interactions and covariates that might

479    be of any interest, Crawley 2013) or as the *global model* (a highly parameterized model

480    containing the variables and associated parameters thought to be important of the

481    problem at hand, Burnham & Anderson 2002; Grueber et al 2011). We shall use the

482    latter term here for consistency with terminology used in information-theory (Grueber et

483    al 2011).

484         Deciding which terms to include in the model requires careful and rigorous *a*

485    *priori* consideration of the system under study. This may appear obvious; however

486    diverse authors have noticed a lack of careful thinking when selecting variables for

487    inclusion in a model (Peters 1991, Chatfield 1995, Burnham & Anderson 2002). Lack of

488    *a priori* consideration, of what models represent, distinguishes rigorous hypothesis

489    testing from 'fishing expeditions' that seek significant predictors among a large group of

490    contenders. Ideally, the global model should be carefully constructed using the

491    researchers' knowledge and understanding of the system such that only predictors likely

492    to be pertinent to the problem at hand are included, rather than including all the

493    predictors the researcher has collected and/or has available. This is a pertinent issue in

494    the age of 'big data', where researchers are often overwhelmed with predictors and risk

495    skipping the important step of *a priori* hypothesis design. In practice, for peer reviewers

496    it is easy to distinguish fishing expeditions from *a priori* hypothesis sets based on the

497    evidence base presented in introductory sections of research outputs.

498

499    **How Complex Should My Global Model Be?**

500         The complexity of the global model will likely be a trade-off between the number

501    of measured observations (the *n* of the study) and the proposed hypotheses about how

502    the measured variables affect the outcome (response) variable. Lack of careful

503    consideration of the parameters to be estimated can result in overparameterised

504  models, where there are insufficient data to estimate coefficients robustly (Southwood &
505  Henderson 2000, Quinn & Keough 2002, Crawley 2013). In simple GLMs,
506  overparameterisation results in a rapid decline in (or absence of) degrees of freedom
507  with which to estimate residual error. Detection of overparameterisation in LMMs can be
508  more difficult because each random effect uses only a single degree of freedom,
509  however the estimation of variance among small numbers of groups can be numerically
510  unstable. Unfortunately, it is common practice to fit a global model that is simply as
511  complex as possible, irrespective of what that model actually represents; that is a
512  dataset containing $k$ predictors yields a model containing a $k$-way interaction among all
513  predictors and simplify from there (Crawley 2013). This approach is flawed for two
514  reasons. First, this practice encourages fitting biologically-unfeasible models containing
515  nonsensical interactions. It should be possible to draw and/or visualise what the fitted
516  model 'looks like' for various combinations of predictors – generally extremely difficult
517  when more than two terms are interacting. Second, using this approach makes it very
518  easy to fit a model too complex for the data. At best, the model will fail to converge, thus
519  preventing inference. At worst, the model will "work", risking false inference. Guidelines
520  for the ideal ratio of data points ($n$) to estimated parameters ($k$) vary widely (see
521  Forstmeier & Schielzeth 2011). Crawley (2013) suggests a minimum $n/k$ of 3, though we
522  argue this is very low and that an n/k of 10 is more conservative. A 'simple' model
523  containing a 3-way interaction between continuous predictors, all that interaction's
524  daughter terms, and a single random intercept needs to estimate 8 parameters, so
525  requires a dataset of a *minimum n* of 80 using this rule. Interactions can be especially
526  demanding, as fitting interactions between a multi-level factor and a continuous
527  predictor can result in poor sample sizes for specific treatment combinations even if the
528  total $n$ is quite large (Zuur et al 2010), which will lead to unreliable model estimates.
529    *Grueber et al (2011) show an excellent worked example of a case where the*
530  *most complex model is biologically feasible and well-reasoned, containing only one 2-*
531  *way interaction. Nakagawa and Foster (2004) discuss the use of power analyses, which*
532  *will be useful in determining the appropriate n/k ratio for a given system.*
533
534  *Assessing Predictor Collinearity*

535    With the desired set of predictors identified, it is wise to check for collinearity among

536    predictor variables. Collinearity among predictors can cause several problems in model

537    interpretation because those predictors explain some of the same variance in the

538    response variable, and their effects cannot be estimated independently (Quinn and

539    Keough. 2002; Graham 2003): First, it can cause model convergence issues as models

540    struggle to partition variance between predictor variables. Second, positively correlated

541    variables can have negatively correlated regression coefficients, as the marginal effect

542    of one is estimated, given the effect of the other, leading to incorrect interpretations of

543    the direction of effects (Figure 2). Third, collinearity can inflate standard errors of

544    coefficient estimates and make 'true' effects harder to detect (Zuur et al 2010). Finally,

545    collinearity can affect the accuracy of model averaged parameter estimates during

546    multi-model inference (Freckleton 2011; Cade 2015). Examples of collinear variables

547    include climatic data such as temperature and rainfall, and morphometric data such as

548    body length and mass. Collinearity can be detected in several ways, including creating

549    correlation matrices between raw explanatory variables, with values >0.7 suggesting

550    both should not be used in the same model (Dormann et al. 2013); or calculating the

551    variance inflation factor (VIF) of each predictor that is a candidate for inclusion in a

552    model (details in Zuur et al 2010) and dropping variables with a VIF higher than a

553    certain value (e.g. 3; Zuur et al 2010, or 10, Quinn & Keogh 2002). One problem with

554    these methods though is that they rely on a user-selected choice of threshold of either

555    the correlation coefficient or the VIF, and use of more stringent (lower) is probably

556    sensible. Some argue that one should always prefer inspection of VIF values over

557    correlation coefficients of raw predictors because strong multicollinearity can be hard to

558    detect with the latter. When collinearity is detected, researchers can either select one

559    variable as representative of multiple collinear variables (Austin 2002), ideally using

560    biological knowledge/ reasoning to select the most meaningful variable (Zuur et al

561    2010); or conduct a dimension-reduction analysis (e.g. Principal Components Analysis;

562    James & McCullugh 1990), leaving a single variable that accounts for most of the

563    shared variance among the correlated variables. Both approaches will only be

564    applicable if it is possible to group explanatory variables by common features, thereby

565    effectively creating broader, but still meaningful explanatory categories. For instance, by

566  using mass and body length metrics to create a 'scaled mass index' representative of

567  body size (Peig & Green 2009).

568

569  *Standardising and Centering Predictors*

570  Transformations of predictor variables are common, and can improve model

571  performance and interpretability (Gelman & Hill 2007). Two common transformations for

572  continuous predictors are i) predictor centering, the mean of predictor $x$ is subtracted

573  from every value in $x$, giving a variable with mean 0 and SD on the original scale of x;

574  and ii) predictor standardising, where $x$ is centred and then divided by the SD of x,

575  giving a variable with mean 0 and SD 1. Rescaling the mean of predictors containing

576  large values (e.g. rainfall measured in thousands of mm) through

577  centering/standardising will often solve convergence problems, in part because the

578  estimation of intercepts is brought into the main body of the data themselves. Both

579  approaches also remove the correlation between main effects and their interactions,

580  making main effects more easily interpretable when models also contain interactions

581  (Schielzeth 2010). Note that this collinearity among coefficients is distinct from

582  collinearity between two separate predictors (see above). Centering and standardising

583  by the mean of a variable changes the interpretation of the model intercept to the value

584  of the outcome expected when $x$ is at its mean value. Standardising further adjusts the

585  interpretation of the coefficient (slope) for $x$ in the model to the change in the outcome

586  variable for a 1 SD change in the value of x. Scaling is therefore a useful tool to improve

587  the stability of models and likelihood of model convergence, and the accuracy of

588  parameter estimates *if* variables in a model are on large (e.g. thousands of mm of

589  rainfall), or vastly different scales. When using scaling, care must be taken in the

590  interpretation and graphical representation of outcomes.

591        *Further reading: Schielzeth (2010) provides an excellent reference to the*

592  *advantages of centering and standardising predictors. Gelman (2008) provides strong*

593  *arguments for standardising continuous variables by 2 SDs when binary predictors are*

594  *in the model. Gelman & Hill (2007 p. 56, 434) discuss the utility of centering by values*

595  *other than the mean.*

596

**Quantifying GLMM Fit and Performance**

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

*Inspection of Residuals and Linear Model Assumptions*

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

628    line (e.g. Crawley 2013), so should a QQ plot of the random effect means (Schielzeth &

629    Nakagawa 2013).

630    *Further reading: Zuur et al (2010) give an excellent overview of the assumptions of*

631    *linear models and how to test for their violation. See also Gelman & Hill (2007 p. 45).*

632    *The R package 'sjPlot' (Lüdecke 2017) has built in functions for several LMM*

633    *diagnostics, including random effect QQ plots. Zuur et al (2009) provides a vast*

634    *selection of model diagnostic techniques for a host of model types, including*

635    *Generalised Least Squared (GLS), GLMMs and Generalized Additive Mixed Effects*

636    *Models (GAMMS).*

637

638    *Overdispersion*

639    Models with a Gaussian (Normal) error structure do not require adjustment for

640    overdispersion, as Gaussian models do not assume a specific mean-variance

641    relationship. For generalized mixed models (GLMMs) however (e.g. Poisson, Binomial),

642    the variance of the data can be greater than predicted by the error structure of the

643    model (e.g. Hilbe 2011). Overdispersion can be caused by several processes

644    influencing data, including zero-inflation, aggregation (non-independence) among

645    counts, or both (Zuur et al 2009). The presence of overdispersion in a model suggests it

646    is a bad fit, and standard errors of estimates will likely be biased unless overdispersion

647    is accounted for (e.g. Harrison 2014). The use of canonical binomial and Poisson error

648    structures, when residuals are overdispersed, tends to result in Type I errors because

649    standard errors are underestimated. Adding an observation-level random effect (OLRE)

650    to overdispersed Poisson or Binomial models can model the overdispersion and give

651    more accurate estimates of standard errors (Harrison 2014; 2015). However, OLRE

652    models may yield inferior fit and/or biased parameter estimates compared to models

653    using compound probability distributions such as the Negative-Binomial for count data

654    (Hilbe 2011; Harrison 2014) or Beta-Binomial for proportion data (Harrison 2015), and

655    so it is good practice to assess the relative fit of both types of model using AIC before

656    proceeding (e.g. Zuur et al 2009). Researchers very rarely report the overdispersion

657    statistic (but see Elston et al 2001), and it should be made a matter of routine. See

658    'Assessing Model Fit Through Simulation' Section for advice on how to quantify and

659    model overdispersion. Note that models can also be underdispersed (*less* variance than
660    expected/predicted by the model, but the tools for dealing with underdispersion are less
661    well developed (Zuur et al 2009). The *spaMM* package (Rousset & Ferdy 2014) can fit
662    models that can handle both overdispersion and underdispersion.
663        *Further reading: Crawley (2013 page 580-581) gives an elegant demonstration of*
664    *how failing to account for overdispersion leads to artificially small standard errors and*
665    *spurious significance of variables. Harrison (2014) quantifies the ability of OLRE to cope*
666    *with overdispersion in Poisson models. Harrison (2015) compares Beta-Binomial and*
667    *OLRE models for overdispersed proportion data.*
668
669    *$R^2$*
670    In a linear modelling context, $R^2$ gives a measure of the proportion of explained variance
671    in the model, and is an intuitive metric for assessing model fit. Unfortunately, the issue
672    of calculating $R^2$ for (G)LMMs is particularly contentious; whereas residual variance can
673    easily be estimated for a simple linear model with no random effects and a Normal error
674    structure, this is not the case for (G)LMMS. In fact, two issues exist with generalising $R^2$
675    measures to (G)LMMs: i) for generalised models containing non-Normal error
676    structures, it is not clear how to calculate the residual variance term on which the $R^2$
677    term is dependent; and ii) for mixed effects models, which are hierarchical in nature and
678    contain error (unexplained variance) at each of these levels, it is uncertain which level to
679    use to calculate a residual error term (Nakagawa & Schielzeth 2013). Diverse methods
680    have been proposed to account for this in GLMMs, including multiple so-called 'pseudo-
681    $r^2$' measures of explained variance (e.g. Nagelkerke 1991, Cox & Snell 1989), but their
682    performance is often unstable for mixed models and can return negative values
683    (Nakagawa & Schielzeth 2013). Gelman & Pardoe (2006) derived a measure of $R^2$ that
684    accounts for the hierarchical nature of LMMs and gives a measure for both group and
685    unit level regressions (see also Gelman & Hill 2007 p. 474), but it was developed for a
686    Bayesian framework and a frequentist analogue does not appear to be widely
687    implemented. The method that has gained the most support over recent years is that of
688    Nakagawa & Schielzeth (2013).

689   The strength of the Nakagawa & Schielzeth (2013) method for GLMMs is that it

690   returns two complementary $R^2$ values: the marginal $R^2$ encompassing variance

691   explained by only the fixed effects, and the conditional $R^2$ comprising variance

692   explained by both fixed and random effects i.e. the variance explained by the whole

693   model (Nakagawa & Schielzeth 2013). Ideally, both should be reported in publications

694   as they provide different information; which one is more 'useful' may depend on the

695   rationale for specifying random effects in the first instance. Recently, Nakagawa,

696   Johnson & Schielzeth (2017) expanded their $R^2$ method to handle models with

697   compound probability distributions like the Negative Binomial error family. Note that

698   when observation-level random effects are included (see 'Overdispersion' section

699   above), the conditional $R^2$ becomes less useful as a measure of explained variance

700   because it includes the extra-parametric dispersion being modelled, but has no

701   predictive power (Harrison 2014).

702   *Further reading: Nakagawa & Schielzeth (2013) provide an excellent and*

703   *accessible description of the problems with, and solutions to, generalising $R^2$ metrics to*

704   *GLMMs. The Nakagawa & Schielzeth (2013) $R^2$ functions have been incorporated into*

705   *several packages, including 'MuMIn' (Bartoń 2016) and 'piecewiseSEM' (Lefcheck*

706   *2015), and Johnson (2014) has developed an extension of the functions for random*

707   *slope models. See Harrison (2014) for a cautionary tale of how the GLMM $R^2$ functions*

708   *are artificially inflated for overdispersed models.*

709

710

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

712   When models are too complex relative to the amount of data available, GLMM variance

713   estimates can collapse to zero (they cannot be negative, not to be confused with

714   *co*variance estimates which can be negative). This is not a problem *per se*, but it's

715   important to acknowledge that in this case the model is equivalent to a standard GLM.

716   Reducing model complexity by removing interactions will often allow random effects

717   variance component estimates to become >0, but this is problematic if quantifying the

718   interaction is the primary goal of the study. REML (restricted/residual maximum

719   likelihood) should be used for estimating variance components of random effects in

720  Gaussian models as it produces less biased estimates compared to ML (maximum
721  likelihood) (Bolker et al 2009). However, when comparing two models with the same
722  random structure but different fixed effects, ML estimation cannot easily be avoided.
723  The RLRsim package (Scheipl, 2016) can be used to calculate restricted likelihood ratio
724  tests for variance components in mixed and additive models. Crucially, when testing the
725  significance of a variance component we are 'testing on the boundary' (Bolker et al
726  2009). That is the null hypothesis for random effects ($\sigma=0$) is at the boundary of its
727  possible range (it has to be $\geq 0$), meaning p-values from a likelihood ratio test are
728  inaccurate. Dividing p values by 2 for tests of single variance components provides an
729  approximation to remedy this problem (Verbenke & Molenberghs, 2000).
730      Finally, estimating degrees of freedom for tests of random effects is difficult, as a
731  random effect can theoretically use anywhere between 1 and $N - 1$ df (where N is the
732  number of random-effect levels) (Bolker et al. 2009). Adequate F and *P* values can be
733  calculated using Satterthwaite or Kenward-Roger approximations to determine
734  denominator degrees of freedom implemented in the package 'lmerTest' (Kuznetzova et
735  al. 2014, see further details in section 'Model Selection and Multi-Model Inference'
736  below).

737

738  *Assessing Model Fit through Simulation*
739  Simulation is a powerful tool for assessing model fit (Gelman & Hill 2007; Kery 2010;
740  Zuur & Ieno 2016), but is rarely used. The premise here is simple: when simulating a
741  response variable from a given set of parameter estimates (a model), the fit of the
742  model to those *simulated* 'ideal' response data should be comparable to the model's fit
743  to the real response variable (Kery 2010). Each iteration yields a simulated dataset that
744  allows calculation of a statistic of interest such as the sum of squared residuals (Kery
745  2010), the overdispersion statistic (Harrison 2014) or the percentage of zeroes for a
746  Poisson model (Zuur & Ieno 2016). If the model is a good fit, after a sufficiently large
747  number of iterations (e.g. 10,000) the distribution of this statistic should encompass the
748  observed statistic in the real data. Significant deviations outside of that distribution
749  indicate the model is a poor fit (Kery 2010). Figure 3 shows an example of using
750  simulation to assess the fit of a Poisson GLMM. After fitting a GLMM to count data, we

751   may wish to check for overdispersion and/or zero-inflation, the presence of which might

752   suggest we need to adjust our modelling strategy. Simulating 10,000 datasets from our

753   model reveals that the proportion of zeroes in our real data is comparable to simulated

754   expectation (Figure 3A). Conversely, simulating 1000 datasets and refitting our model to

755   each dataset, we see that the sum of the squared Pearson residuals for the real data is

756   far larger than simulated expectation (Figure 3B), giving evidence of overdispersion

757   (Harrison 2014). We can use the simulated frequency distribution of this test statistic to

758   derive a mean and 95% confidence interval for the overdispersion by calculating the

759   ratio of our test statistic to the simulated values (Harrison 2014). The dispersion statistic

760   for our model is 3.16 [95% CI 2.77 – 3.59]. Thus, simulations have allowed us to

761   conclude that our model is overdispersed, but that this overdispersion is not due to

762   zero-inflation. All R code for reproducing these simulations is provided in Online

763   Supplementary Material.

764         *Further reading: The R package 'SQuiD' (Allegue et al 2017) provides a highly*

765   *flexible simulation tool for learning about, and exploring the performance of, GLMMs.*

766   *Rykiel (1996) discusses the need for validation of models in ecology.*

767

768   *Dealing with missing data*

769   When collecting ecological data it is often not possible to measure all of the predictors

770   of interest for every measurement of the dependant variable. Such missing data are a

771   common feature of ecological datasets, however the impacts of this have seldom been

772   considered in the literature (Nakagawa & Freckleton 2011). Incomplete rows of data in

773   dataframes i.e. those missing predictor and/or response variables are often dealt with

774   by deleting or ignoring those rows of data entirely when modelling (Nakagawa &

775   Freckleton 2008), although this may result in biased parameter estimates and,

776   depending on the mechanism underlying the missing data, reduces statistical power

777   (Nakagawa & Freckleton 2008). Nakagawa & Freckleton (2011) recommend multiple

778   imputation (MI) as a mechanism for handling non-informative missing data, and

779   highlight the ability of this technique for more accurate estimates, particularly for

780   information theoretic / AIC approaches.

781   *Further reading: See Nakagawa & Freckleton (2008) for a review on the risks of*

782   *ignoring incomplete data. Nakagawa & Freckleton (2011) demonstrate the effects of*

783   *missing data during model selection procedures, and provide an overview of R*

784   *packages available for MI. Nakagawa (2015)* and Noble & Nakagawa (2017) discuss

785   methods for dealing with missing data in ecological statistics.

786   ## Model Selection and Multi-Model Inference

787   Model selection seeks to optimise the trade-off between the fit of a model given the data

788   and that model's complexity. Given that the researcher has a robust global model that

789   satisfies standard assumptions of error structure and hierarchical independence,

790   several methods of model selection are available, each of which maximises the fit-

791   complexity trade off in a different way (Johnson & Omland 2004). We discuss the

792   relative merits of each approach briefly here, before expanding on the use of

793   information-theory and multi-model inference in ecology. We note that these

794   discussions are not meant to be exhaustive comparisons, and we encourage the reader

795   to delve into the references provided for a comprehensive picture of the arguments for

796   and against each approach.

797

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

799   A common approach to model selection is the comparison of a candidate model

800   containing a term of interest to the corresponding 'null' model lacking that term, using a

801   p value from a likelihood ratio test (LRT), referred to as null-hypothesis significance

802   testing (NHST; Nickerson 2000). Stepwise deletion is a model selection technique that

803   drops terms sequentially from the global model to  arrive at a 'minimal adequate model'

804   (MAM). Evaluating whether a term should be dropped or not can be done using NHST

805   to arrive at a model containing only significant predictors (see Crawley 2013), or using

806   information theory to yield a model containing only terms that cause large increases in

807   information criterion score when removed. Stepwise selection using NHST is by far the

808   most common variant of this approach, and so we focus on this method here.

809   Stepwise deletion procedures have come under heavy criticism; they can overestimate

810   the effect size of significant predictors (Whittingham et al 2006; Forstmeier & Schielzeth

811    2011; Burnham, Anderson & Huyvaert 2011) and force the researcher to focus on a

812    single best model as if it were the only combination of predictors with support in the

813    data. Although we strive for simplicity and parsimony, this assumption is not always

814    reasonable in complex ecological systems (e.g. Burnham, Anderson & Huyvaert 2011).

815    It is common to present the MAM as if it arose from a single *a priori* hypothesis, when in

816    fact arriving at the MAM required multiple significance tests (Whittingham et al 2006;

817    Forstmeier & Schielzeth 2011). This cryptic multiple testing can lead to hugely inflated

818    Type I errors (Forstmeier & Schielzeth 2011). Perhaps most importantly, LRT can be

819    unreliable for fixed effects in GLMMs unless both total sample size and replication of the

820    random effect terms is high (see Bolker et al 2009 and references therein), conditions

821    which are often not satisfied for most ecological datasets. Because stepwise deletion

822    can cause biased effect sizes, presenting means and SEs of parameters from the global

823    model should be more robust, especially when the n/k ratio is low (Forstmeier &

824    Schielzeth 2011). Performing 'full model tests' (comparing the global model to an

825    intercept only model) before investigating single-predictor effects controls the Type I

826    error rate (Forstmeier & Schielzeth 2011). Reporting the full model also helps reduce

827    publication bias towards strong effects, providing future meta-analyses with estimates of

828    both significant and non-significant effects (Forstmeier & Schielzeth 2011). Global

829    model reporting should not replace other model selection methods, but provides a

830    robust measure of how likely significant effects are to arise by sampling variation alone.

831         *Further reading: See Murtaugh's (2014) excellent 'in Defense of P values', as*

832    *well as the other papers on the topic in the same special issue of Ecology. Stephens et*

833    *al (2005) & Mundry (2011) argue the case for NHST under certain circumstances such*

834    *as well-designed experiments. Halsey et al (2015) discuss the wider issues of the*

835    *reliability of p values relative to sample size.*

836

837    *Information-Theory and Multi-Model Inference*

838    Unlike NHST, which leads to a focus on a single best model, model selection using

839    information theoretic (IT) approaches allows the degree of support in the data for

840    several competing models to be ranked using metrics such as Akaike's Information

841    Criterion (AIC). Information criteria attempt to quantify the Kullback-Leibler distance

842   (KLD), a measure of the relative amount of information lost when a given model

843   approximates the true data-generating process. Thus, relative difference among models

844   in AIC should be representative in relative differences in KLD, and the model with the

845   lowest AIC should lose the least information and be the best model in that it optimises

846   the trade-off between fit and complexity (e.g. Richards 2008). A key strength of the IT

847   approach is that it accounts for 'model selection uncertainty', the idea that several

848   competing models may all fit the data similarly well (Burnham & Anderson 2002;

849   Burnham, Anderson & Huyvaert 2011). This is particularly useful when competing

850   models share equal "complexity" (i.e. number of predictors, or number of residual

851   degrees of freedom): in such situations, NHST is impossible because NHST requires a

852   simpler (nested) model for comparison. Where several models have similar support in

853   the data, inference can be made from all models using model-averaging (Burnham &

854   Anderson 2002; Johnson & Omand 2004; Grueber et al 2011). Model averaging

855   incorporates uncertainty by weighting the parameter estimate of a model by that

856   model's Akaike weight (often referred to as the probability of that model being the best

857   Kullback-Leibler model given the data, but see Richards 2005). Multi-model inference

858   places a strong emphasis on *a priori* formulation of hypotheses (Burnham & Anderson

859   2002; Dochterman & Jenkins 2011; Lindberg et al 2015), and model-averaged

860   parameter estimates arising from multi-model inference are thought to lead to more

861   robust conclusions about the biological systems compared to NHST (Johnson &

862   Omland 2004, but see Richards et al 2011). These strengths over NHST have meant

863   that the use of IT approaches in ecology and evolution has grown rapidly in recent years

864   (Lindberg et al 2015; Barker & Link 2015; Cade 2015). We do not expand on the

865   specific details of the difference between NHST and IT here, but point the reader to

866   some excellent references on the topic. Instead, we use this section to highlight recent

867   empirical developments in the best practice methods for the application of IT in ecology

868   and evolution.

869          *Further reading: Grueber et al (2011) and Symonds & Moussalli (2011) give a*

870   *broad overview of multi-model inference in ecology, and provide a worked model*

871   *selection exercise. Heygi & Garamszegi (2011) provide a detailed comparison of IT and*

872   *NHST approaches. Burnham, Anderson & Huyvaert (2011) demonstrate how AIC*

873 *approximates Kullback-Leibler information and provide some excellent guides for the*

874 *best practice of applying IT methods to biological datasets. Vaida & Blanchard (2005)*

875 *provide details on how AIC should be implemented for the analysis of clustered data.*

876

877

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

879

880     *1.  Using All-Subsets Selection*

881 All-Subsets selection is the act of fitting a global model, often containing every possible

882 interaction, and then fitting every possible nested model. On the surface, all-subsets

883 might appear to be a convenient and fast way of 'uncovering' the causal relationships in

884 the data. All-subsets selection of enormous global models containing large numbers of

885 predictors and their interactions makes analyses extremely prone to including

886 uninformative parameters and 'overfitted' models. Burnham & Anderson (2002) caution

887 strongly against all-subsets selection, and instead advocate 'hard thinking' about the

888 hypotheses underlying the data. If adopting an all subsets approach, it is worth noting

889 the number of models to consider increases exponentially with the number of predictors,

890 where 5 predictors require $2^5$ (32) models to be fitted, whilst 10 predictors requires 1024

891 models, both *without* including any interactions but including the null model.

892       Global models should not contain huge numbers of variables and interactions

893 without prior thought about what the models represent for a study system. In cases

894 where all-subsets selection from a global model is performed, it is important to view

895 these model selection exercises as exploratory (Symonds & Moussali 2011), and hold

896 some data back from these exploratory analyses to be used for cross-validation with the

897 top model(s) (see Dochterman and Jenkins 2011 and references therein). Here, 90% of

898 the data can be used to fit the model(s), with the remaining 10% used for confirmatory

899 analysis to quantify how well the model(s) perform for prediction (Zuur & Ieno 2016).

900 Such an approach requires a huge amount of data (Dochterman and Jenkins 2011), but

901 cross-validation to validate a model's predictive ability is rare and should result in more

902 robust inference (see also Fieberg & Johnson 2015).

903    Therefore, best practice is to consider only a handful of hypotheses and then build a

904    single statistical model to reflect each hypothesis. This makes inference easier because

905    the resulting top model set will likely contain fewer parameters, and certainly fewer

906    uninformative parameters (Burnham & Anderson 2002; Arnold 2010). However, we

907    argue all subsets selection may be sensible in a limited number of circumstances when

908    testing causal relationships between explanatory variables and the response variable.

909    For example, if the most complex model contains two main effects and their interaction,

910    performing all subsets selection on that model is identical to building the five competing

911    models (including the null model) nested in the global model, all of which may be

912    considered likely to be supported by the data. A small number of models built to reflect

913    well-reasoned hypotheses are only valid if the predictors therein are not collinear (see

914    'Collinearity' section above). All-subsets selection using the R package *MuMIn (*Bartoń

915    2016*)* will not automatically check for collinearity, and so the onus falls on the

916    researcher to be thorough in checking for such problems.

917

918    *2.  Deciding Which Information Criterion To Use*

919    Several information criteria are available to rank competing models, but their

920    calculations differ subtly. Commonly applied criteria include Akaike's Information

921    Criterion (AIC), the small sample size correction of AIC for when n/k <40 (AICc), and the

922    Bayesian Information Criterion (BIC). QAIC is an adjustment to AIC that accounts for

923    overdispersion, and should be used when overdispersion has been identified in a model

924    (see 'Overdispersion section' above). Note QAIC is not required if the overdispersion in

925    the dataset has been modelled using zero-inflated models, observation-level random

926    effects, or compound probability distributions. Bolker et al (2009) and Grueber et al

927    (2011) provide details of how to calculate these criteria.

928       AIC maximises the fit/complexity trade-off of a model by balancing the model fit

929    with the number of estimated parameters. AICc and BIC both penalise the IC score

930    based on total sample size *n*, but the degree of penalty for AICc is less severe than BIC

931    for moderate sample sizes, and more severe for very low sample size (Brewer et al

932    2016). Whilst AIC tend to select overly complex models, Burnham and Anderson (2002)

933    criticised BIC for selecting overly simplistic models (underfitting). BIC is also criticised

934    because it operates on the assumption that the true model is in the model set under
935    consideration, whereas in ecological studies this is unlikely to be true (Burnham &
936    Anderson 2002; 2004). Issues exist with both AIC and BIC in a GLMM context for
937    estimating the number of parameters for a random effect (Bolker et al 2009; Grueber et
938    al 2011), and although degrees of freedom corrections to remedy this problem exist it is
939    not always clear what method is being employed by software packages (see Bolker et al
940    2009 Box 3). Brewer et al (2016) show how the optimality of AIC, AICc and BIC for
941    prediction changes with both sample size and effect size of predictors (see also
942    Burnham and Anderson 2004). Therefore, the choice between the two metrics is not
943    straightforward, and may depend on the goal of the study i.e. model selection vs
944    prediction, see Grueber et al 2011 Box 1.
945
946        *3.  Choice of ΔAIC Threshold*
947    Once all models have been ranked by an information criterion, it is common practice to
948    identify a "top model set" containing all models assumed to have comparable support in
949    the data, normally based on the change in AIC values relative to the best AIC model
950    (ΔAIC). Historically, Burnham & Anderson (2002) recommended that only models with
951    ΔAIC between 0-2 should be used for inference, but subsequent work has shown that
952    for some models a much higher ΔAIC  cut off is required to give a 95% probability of
953    including the best (expected) Kullback-Leibler Distance model in the top model set
954    (Richards 2008; see also Burnham et al 2011). An alternative approach to using ΔAIC
955    cut offs is to include all models within a cumulative Akaike weight  of ≥0.95 from the top
956    model in the "95% confidence set" (Burnham & Anderson 2002; Symonds & Moussali
957    2011). Using high cut-offs is not encouraged, to avoid overly complex model sets
958    containing uninformative predictors (Richards 2008; Grueber et al. 2011) but deciding
959    on how many is too many remains a contentious issue (Grueber et al. 2011). We
960    suggest Δ6 as a minimum following Richards (2005; 2008).
961
962        *4.  Using the Nesting Rule to Improve Inference from the Top Model Set*
963    It is well known that AIC tends towards overly complex models ('overfitting', Burnham &
964    Anderson 2002). As AIC only adds a 2 point penalty to a model for inclusion of a new

965    term, Arnold (2010) demonstrated that adding a nuisance (completely random) predictor

966    to a well-fitting model leads to a ΔAIC value of the new model of ~ 2, therefore

967    appearing to warrant inclusion in the top model set (see section above). Therefore,

968    inference can be greatly improved by eliminating models from the top model set that are

969    more complex versions of nested models with better AIC support, known as the nesting

970    rule (Richards 2005; 2008; Richards et al2011). Doing so greatly reduces the number of

971    models to be used for inference, and improves parameter accuracy (Arnold 2010;

972    Richards et al 2008). Symonds & Moussali (2011) caution that its applicability has not

973    yet been widely assessed over a range of circumstances, but the theory behind its

974    application is sound and intuitive (Arnold 2010).

975

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

977    With a top model set in hand, it is common practice to use the summed Akaike weights

978    of every model in that set in which a predictor of interest occurs as a measure of

979    'variable importance' (e.g. Grueber et al 2011). Recent work has demonstrated that this

980    approach is flawed because Akaike weights are interpreted as relative model

981    probabilities, and give no information about the importance of individual predictors in a

982    model (Cade 2015), and fail to distinguish between variables with weak or strong effects

983    (Galipaud et al 2014; 2017). The sum of Akaike weights as a measure of variable

984    importance may at best be a measure of how likely a variable would be included in the

985    top model set after repeated sampling of the data (Burnham & Anderson 2002; Cade

986    2015, but see Galipaud et al 2017). A better measure of variable importance would be

987    to compare standardised effect sizes (Schielzeth 2010; Cade 2015). However, summed

988    Akaike weights for variables in top model sets still represent useful quantitative

989    evidence (Giam & Olden 2016); they should be reported in model summary tables, and

990    ideally interpreted in tandem with model averaged effect sizes for individual parameters.

991

992        *6.  Model Averaging when Predictors Are Collinear*

993    The aim of model averaging is to incorporate the uncertainty in the size and presence of

994    effects among a set of candidate models with similar support in the data. Model

995    averaging using Akaike weights proceeds on the assumption that predictors are on

996  common scales across models and are therefore comparable. Unfortunately, the nature

997  of multiple regression means that the scale and sign of coefficients will change across

998  models depending on the presence or absence of other variables in a focal model

999  (Cade 2015). The issue of predictor scaling changing across models is particularly

1000  exacerbated when predictors are collinear, even when VIF values are low (Burnham

1001  and Anderson 2002; Lukacs, Burnham & Anderson 2010; Cade 2015). Cade (2015)

1002  recommends standardising model parameters based on partial standard deviations to

1003  ensure predictors are on common scales across models prior to model averaging

1004  (details in Cade 2015). We stress again the need to assess multicollinearity among

1005  predictors in multiple regression modelling before fitting models (Zuur & Ieno 2016) and

1006  before model-averaging coefficients from those models (Lukacs, Burnham & Anderson

1007  2010; Cade 2015)

1008

1009

# 1010  Conclusion

1011  We hope this article will act as both a guide, and as a gateway to further reading, for

1012  both new researchers and those wishing to update their portfolio of analytic techniques.

1013  Here we distil our message into a bulleted list.

1014  1. Modern mixed effect models offer an unprecedented opportunity to explore complex

1015  biological problems by explicitly modelling non-Normal data structures and/or non-

1016  independence among observational units. However, the LMM and GLMM toolset should

1017  be used with caution.

1018  2. Rigorous testing of both model fit ($R^2$) and model adequacy (violation of assumptions

1019  like homogeneity of variance) must be carried out. We must recognise that satisfactory

1020  fit does not guarantee we have not violated the assumptions of LMM, and vice versa.

1021  Interpret measures of $R^2$ for (G)LMMs with hierarchical errors cautiously, especially

1022  when OLRE are used.

1023  3. Collinearity among predictors is difficult to deal with and can severely impair model

1024  accuracy. Be especially vigilant if data are from field surveys rather than controlled

1025  experiments, as collinearity is likely to be present.

1026   4. When including a large number of predictors is necessary, backwards selection and

1027   NHST should be avoided, and ranking via AIC of all competing models is preferred. A

1028   critical question that remains to be addressed is whether model selection based on

1029   information theory is superior to NHST even in cases of balanced experimental designs

1030   with few predictors.

1031   5. Data simulation is a powerful but underused tool. If the analyst harbours any

1032   uncertainty regarding the fit or adequacy of the model structure, then the analysis of

1033   data simulated to recreate the perceived structure of the favoured model can provide

1034   reassurance, or justify doubt.

1035   6. Wherever possible, provide diagnostic assessment of model adequacy, and metrics

1036   of model fit, even if in Supplementary Material.

1037

## 1038 Acknowledgements

1041
1042

# References

1043

1044 Aarts E, Dolan CV, Verhage M, Sluis S. 2015. Multilevel analysis quantifies variation in
1045       the experimental effect while optimizing power and preventing false
1046       positives. *BMC Neuroscience* 16:94.

1047 Allegue H, Araya-Ajoy YG, Dingemanse NJ, Dochtermann NA, Garamszegi LZ,
1048       Nakagawa S, Reale D, Schielzeth H, Westneat DF. 2017. Statistical Quantification
1049       of Individual Differences (SQuID): an educational and statistical tool for
1050       understanding multilevel phenotypic data in linear mixed models. *Methods in*
1051       *Ecology and Evolution* 8:257-67.

1052 Arnold TW. 2010. Uninformative parameters and model selection using Akaike's
1053       Information Criterion. The *Journal of Wildlife Management* 74: 1175-1178.

1054 Austin MP. 2002. Spatial prediction of species distribution: an interface between
1055       ecological theory and statistical modelling. *Ecological Modelling* 157: 101–118.

1056 Barker RJ, Link WA. 2015. Truth, models, model sets, AIC, and multimodel inference: A
1057       Bayesian perspective. *The Journal of Wildlife Management* 79: 730–738.

1058 Barr DJ, Levy R, Scheepers C, Tily HJ. 2013. Random effects structure for confirmatory
1059       hypothesis testing: Keep it maximal. *Journal of memory and language* 68:255-78.

1060 Bartoń K. 2016. MuMIn: Multi-Model Inference. R package version
1061       1.15.6.https://CRAN.R-project.org/package=MuMIn

1062 Bates D, Maechler M, Bolker B, Walker S. 2015. Fitting Linear Mixed-Effects Models
1063       Using lme4. *Journal of Statistical Software* 67: 1-48.

1064 Bates D, Kliegl R, Vasishth S, Baayen H. 2015. Parsimonious mixed models. *arXiv*
1065       *preprint arXiv:1506.04967*.

1066 Bolker BM, Brooks ME, Clark CJ, Geange SW, Poulsen JR, Stevens MHH, White JSS.
1067       2009. Generalized linear mixed models: a practical guide for ecology and
1068       evolution. *Trends in Ecology and Evolution* 24: 127–135.

1069 Breslow NE, Clayton DG. 1993. Approximate inference in generalized linear mixed
1070       models. *Journal of the American Statistical Association* 88: 9-25.

1071 Brewer MJ, Butler A, Cooksley SL. 2016. The relative performance of AIC, AICC and
1072       BIC in the presence of unobserved heterogeneity. *Methods in Ecology and*
1073       *Evolution* 7: 679-692.

Burnham KP, Anderson DR. 2002. Model Selection and Multimodel Inference: A
Practical Information-Theoretic Approach, Second. Springer-Verlag, New York.

Burnham KP, Anderson DR. 2004. Multimodel inference: understanding AIC and BIC in
model selection. *Sociological Methods & Research* 33: 261-304.

Burnham KP, Anderson DR, Huyvaert KP. 2011. AIC model selection and multimodel
inference in behavioral ecology: Some background, observations, and
comparisons. *Behavioral Ecology and Sociobiology* 65: 23–35.

Cade BS. 2015. Model averaging and muddled multimodel inferences. *Ecology* 96:
2370–2382.

Chatfield C. 1995. Model uncertainty, data mining and statistical inference (with
discussion). *Journal of the Royal Statistical Society, Series A* 158: 419-66.

Cox DR, Snell EJ. 1989. *The Analysis of Binary Data,* 2nd ed. London: Chapman and
Hall.

Crawley (2013) *The R Book*. Second Edition. Wiley, Chichester UK.

Dochtermann NA, Jenkins SH. 2011. Developing multiple hypotheses in behavioural
ecology. *Behavioral Ecology and Sociobiology* 65: 37-45.

Dominicus A, Skrondal A, Gjessing HK, Pedersen NL, Palmgren J. 2006.Likelihood ratio
tests in behavioral genetics: problems and solutions. *Behavior Genetics* 36: 331–
340.

Dormann CF, Elith J, Bacher S, Buchmann C, Carl G, Carré G, Marquéz JR, Gruber B,
Lafourcade B, Leitão PJ, Münkemüller T. 2013. Collinearity: a review of methods
to deal with it and a simulation study evaluating their performance. *Ecography* 36:
027–046.

Ellison AM. 2004. Bayesian inference in ecology. *Ecology letters* 7: 509-520.

Elston, DA, Moss R, Boulinier T, Arrowsmith C, Lambin X, 2001. Analysis of
aggregation, a worked example: numbers of ticks on red grouse
chicks. *Parasitology* 122: 563-569.

Fieberg J, Johnson DH. 2015. MMI: Multimodel inference or models with management
implications? *The Journal of Wildlife Management* 79: 708–718.

1103  Forstmeier W, Schielzeth H. 2011. Cryptic multiple hypotheses testing in linear models:
1104         Overestimated effect sizes and the winner's curse. *Behavioral Ecology and*
1105         *Sociobiology* 65: 47–55.

1106  Freckleton RP. 2011. Dealing with collinearity in behavioural and ecological data: model
1107         averaging and the problems of measurement error. *Behavioral Ecology and*
1108         *Sociobiology* 65: 91-101.

1109  Galipaud M, Gillingham MAF, David M, Dechaume-Moncharmont FX. 2014. Ecologists
1110         overestimate the importance of predictor variables in model averaging: a plea for
1111         cautious interpretations. *Methods in Ecology and Evolution* 5, 983-991.

1112  Galipaud M, Gillingham MAF, Dechaume-Moncharmont FX. 2017. A farewell to the sum
1113         of Akaike weights: The benefits of alternative metrics for variable importance
1114         estimations in model selection. *Methods in Ecology and Evolution* 00:1–11.
1115         https://doi.org/10.1111/2041-210X.12835

1116  Gelman A, Hill J. 2007. *Data analysis using regression and hierarchical/multilevel*
1117         *models*. New York, NY, USA: Cambridge University Press.

1118  Gelman A. 2008. Scaling regression inputs by dividing by two standard
1119         deviations. *Statistics in Medicine 27*: 2865-2873.

1120  Gelman A, Pardoe I. 2006. Bayesian measures of explained variance and pooling in
1121         multilevel (hierarchical) models. *Technometrics* 48: 241-251.

1122  Giam, X., Olden, J.D. (2016) Quantifying variable importance in a multimodel inference
1123         framework. *Methods in Ecology and Evolution*, 7:388-397.

1124  Graham ME (2003) Confronting multicollinearity in multiple linear regression. *Ecology*
1125         84: 2809-2815

1126  Grueber CE, Nakagawa S, Laws RJ, Jamieson IG. 2011. Multimodel inference in
1127         ecology and evolution: Challenges and solutions. *Journal of Evolutionary Biology*
1128         24: 699–711.

1129  Harrison XA. 2014. Using observation-level random effects to model overdispersion in
1130         count data in ecology and evolution. *PeerJ* 2: e616.

1131  Harrison XA. 2015. A comparison of observation-level random effect and Beta-Binomial
1132         models for modelling overdispersion in Binomial data in ecology &
1133         evolution. *PeerJ, 3*: p.e1114.

1134  Halsey LG, Curran-Everett D, Vowler SL, Drummond GB. 2015. The fickle P value
1135      generates irreproducible results. *Nature Methods* 12: 179-185.
1136  Hegyi G, Garamszegi LZ. 2011. Using information theory as a substitute for stepwise
1137      regression in ecology and behaviour. *Behavioral Ecology and Sociobiology* 65: 69-
1138      76.
1139  Hilbe JM. 2011. *Negative binomial regression*. Cambridge University Press.
1140  Houslay T, Wilson A. 2017. Avoiding the misuse of BLUP in behavioral ecology.
1141      *Behavioral Ecology* arx023 doi:10.1093/beheco/arx023
1142  Ives AR. 2015. For testing the significance of regression coefficients, go ahead and
1143      log-transform count data. *Methods in Ecology and Evolution* 6:, 828-835.
1144  James FC, McCullugh CF. 1990. Multivariate Analysis In Ecology And Systematics:
1145      Panacea Or Pandora Box. *Annual Review of Ecology and Systematics* 21: 129–
1146      166.
1147  Johnson JB, Omland KS. 2004. Model selection in ecology and evolution. *Trends in*
1148      *Ecology and Evolution* 19: 101–108.
1149  Johnson PCD. 2014. Extension of Nakagawa & Schielzeth's $R^2$ GLMM to random
1150      slopes models. *Methods in Ecology and Evolution* 5:  944-946.
1151  Kass RE, Caffo BS, Davidian M, Meng XL, Yu B, Reid N. 2016. Ten simple rules for
1152      effective statistical practice. *PLoS Computational Biology* 12: p.e1004961.
1153  Keene ON. 1995. The log transform is special. *Statistics in Medicine* 14: 811–819.
1154  Kéry M. 2010. *Introduction to WinBUGS for ecologists: Bayesian approach to*
1155      *regression, ANOVA, mixed models and related analyses*. Academic Press.
1156  Kuznetsova A, Brockhoff PB, Christensen RHB. 2014. Package 'lmerTest'. Test for
1157      random and fixed effects for linear mixed effect models (lmer objects of lme4
1158      package). R package ver.2.
1159  Lefcheck JS. 2015. piecewiseSEM: Piecewise structural equation modeling in R for
1160      ecology, evolution, and systematics. *Methods in Ecology and Evolution* 7: 573-
1161      579.
1162  Lindberg MS, Schmidt JH, Walker J. 2015. History of multimodel inference via model
1163      selection in wildlife science. *The Journal of Wildlife Management* 79: 704–707.

1164   Low-Décarie E, Chivers C, Granados M. 2014. Rising complexity and falling explanatory
1165         power in ecology. *Frontiers in Ecology and the Environment* 12: 412-418.

1166   Lüdecke D. 2017. SjPlot: Data Visualization for Statistics in Social Science. 2017 *R*
1167         *package version*, 2.4.0.

1168   Lukacs PM, Burnham KP, Anderson DR. 2010. Model selection bias and Freedman's
1169         paradox. *Annals of the Institute of Statistical Mathematics* 62: 117–125.

1170   Mundry R. 2011. Issues in information theory-based statistical inference—a
1171         commentary from a frequentist's perspective. *Behavioral Ecology and*
1172         *Sociobiology* 65: 57-68.

1173   Murtaugh PA. 2007. Simplicity and complexity in ecological data analysis. *Ecology* 88:
1174         56-62.

1175   Murtaugh PA. 2009. Performance of several variable-selection methods applied to real
1176         ecological data. *Ecology Letters* 10: 1061-1068.

1177   Murtaugh PA. 2014. In defense of P values. *Ecology 95*: 611-617

1178   Nagelkerke NJ. 1991. A note on a general definition of the coefficient of determination.
1179         *Biometrika* 78: 691-692.

1180   Nakagawa S, Foster T. 2004. The case against retrospective statistical power analyses
1181         with an introduction to power analysis. *Acta Ethologica* 7: 103-108.

1182   Nakagawa S, Freckleton RP. 2008. Missing inaction: the dangers of ignoring missing
1183         data. *Trends in Ecology and Evolution* 23(11): 592-596.

1184   Nakagawa S, Freckleton RP. 2011. Model averaging, missing data and multiple
1185         imputation: a case study for behavioural ecology. *Behavioral Ecology and*
1186         *Sociobiology* 65: 103-116.

1187   Nakagawa S, Schielzeth H. 2010. Repeatability for Gaussian and non-Gaussian data: a
1188         practical guide for biologists. *Biological Reviews* 85: 935-956

1189   Nakagawa S, Schielzeth H. 2013. A general and simple method for obtaining R2 from
1190         generalized linear mixed-effects models. *Methods in Ecology and Evolution* 4:
1191         133-142.

1192   Nakagawa, S. (2015) *Missing data: mechanisms, methods and messages In: Ecological*
1193         *Statistics: contemporary theory and application* (eds. Fox, G. A., Negrete-
1194         Yankelevich, S. & Sosa, V. J.). Oxford University Press, Oxford. pp. 81-105

1195    Nakagawa S., Johnson PC, Schielzeth H. 2017. The coefficient of determination $R^2$ and

1196          intra-class correlation coefficient from generalized linear mixed-effects models

1197          revisited and expanded. *Journal of The Royal Society Interface 14*(134),

1198          p.20170213.

1199    Nickerson RS. 2000. Null Hypothesis Significance Testing: A Review of an Old and

1200          Continuing Controversy. *Psychological Methods* 5: 241-301.

1201    Noble, D. W. A. & Nakagawa, S. (submitted) Planned missing data design: stronger

1202          inferences increased research efficiency and improved animal welfare in ecology

1203          and evolution. bioRxiv https://www.biorxiv.org/content/early/2018/01/11/247064

1204    O'Hara RB, Kotze DJ. 2010. Do not log-transform count data. *Methods in Ecology and*

1205          *Evolution* 1: 118-122.

1206    Peters RH. 1991. *A critique for ecology*. Cambridge University Press.

1207    Peig J, Green AJ. 2009. New perspectives for estimating body condition from

1208          mass/length data: the scaled mass index as an alternative method. *Oikos* 118:

1209          1883-1891.

1210    Quinn GP, Keough MJ. 2002. *Experimental design and data analysis for biologists*.

1211          Cambridge University Press.

1212    R Core Team. 2016. R: A language and environment for statistical computing. R

1213          Foundation for Statistical Computing, Vienna, Austria. URL https://www.R-

1214          project.org/.

1215    Richards SA. 2005. Testing ecological theory using the information-theoretic approach:

1216          examples and cautionary results. *Ecology* 86: 2805-2814.

1217    Richards SA. 2008. Dealing with overdispersed count data in applied ecology. *Journal*

1218          *of Applied Ecology* 45 218–227.

1219    Richards, SA, Whittingham MJ, Stephens PA. 2011. Model selection and model

1220          averaging in behavioural ecology: the utility of the IT-AIC framework. *Behavioral*

1221          *Ecology and Sociobiology* 65: 77–89.

1222    Rousset F, Ferdy JB 2014. Testing environmental and genetic effects in the presence of

1223          spatial autocorrelation. *Ecography* 37: 781-790

1224    Rykiel EJ. 1996. Testing ecological models: The meaning of validation. *Ecological*

1225          *Modelling* 90: 229-244.

1226  Satterthwaite FE. 1946. An approximate distribution of estimates of variance

1227      components. *Biometrics Bulletin* 2(6): 110-114.

1228  Scheipl F, & Bolker, B. 2016. RLRsim: Exact (Restricted) Likelihood Ratio Tests for

1229      Mixed and Additive Models *Computational Statistics & Data Analysis*. R package

1230      version 3.1-3. https://cran.r-project.org/web/packages/RLRsim/index.html

1231  Schielzeth H, Forstmeier W. 2009. Conclusions beyond support: overconfident

1232      estimates in mixed models. *Behavioral Ecology* 20: 416-420.

1233  Schielzeth H, Nakagawa S. 2013. Nested by design: model fitting and interpretation in a

1234      mixed model era. *Methods in Ecology Evolution* 4: 14-24

1235  Schielzeth H. 2010. Simple means to improve the interpretability of regression

1236      coefficients. *Methods in Ecology and Evolution* 1:  103-113

1237  Southwood TRE, Henderson PA. 2000. *Ecological methods*. John Wiley &

1238      Sons.Stephens PA, Buskirk SW, Hayward GD, Martinez Del Rio C. 2005.

1239      Information theory and hypothesis testing: a call for pluralism. *Journal of Applied*

1240      *Ecology* 42: 4-12.

1241  Symonds MRE, Moussalli A. 2011. A brief guide to model selection, multimodel

1242      inference and model averaging in behavioural ecology using Akaike's information

1243      criterion. *Behavioral Ecology and Sociobiology* 65: 13–21.

1244  Vaida F, Blanchard S. 2005. Conditional Akaike information for   mixed-effects models.

1245      *Biometrika* 92: 351–370

1246  van de Pol M, Wright J. 2009. A simple method for distinguishing within-versus

1247      between-subject effects using mixed models. *Animal Behaviour* 77: 753-758.

1248  Verbenke G, Molenberghs G. 2000. *Linear mixed models for longitudinal data*. New

1249      York, Springer.

1250  Warton D, Hui F. 2011. The arcsine is asinine: the analysis of proportions in ecology.

1251      *Ecology* 92: 3-10

1252  Warton DI, Lyons M, Stoklosa J, Ives AR. 2016. Three points to consider when

1253      choosing a LM or GLM test for count data. *Methods in Ecology and Evolution* 7:

1254      882-90.

1255    Wilson AJ, Réale D, Clements MN, Morrissey MM, Postma E, Walling CA, Kruuk LEB,

1256        Nussey DH. 2010. An ecologist's guide to the animal model. *Journal of Animal*

1257        *Ecology* 79: 13–26.

1258    Wood SN, Goude Y, Shaw S. 2015. Generalized additive models for large data

1259        sets. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 64:139-

1260        155.

1261    Whittingham MJ, Stephens PA, Bradbury RB, Freckleton RP 2006. Why do we still use

1262        stepwise modelling in ecology and behaviour? *Journal of Animal Ecology* 75:

1263        1182-1189.

1264    Zuur AF, Ieno EN, Walker NJ, Saveliev AA, Smith GM. 2009 *Mixed Effects Models and*

1265        *Extensions in Ecology with R* Springer, New York

1266    Zuur AF, Ieno EN, Elphick CS. 2010. A protocol for data exploration to avoid common

1267        statistical problems. *Methods in Ecology and Evolution* 1: 3-14.

1268    Zuur AF, Ieno EN, 2016. A protocol for conducting and presenting results of

1269        regression-type analyses. *Methods in Ecology and Evolution* 7: 636-645.

1270

1271

# Figure 1(on next page)

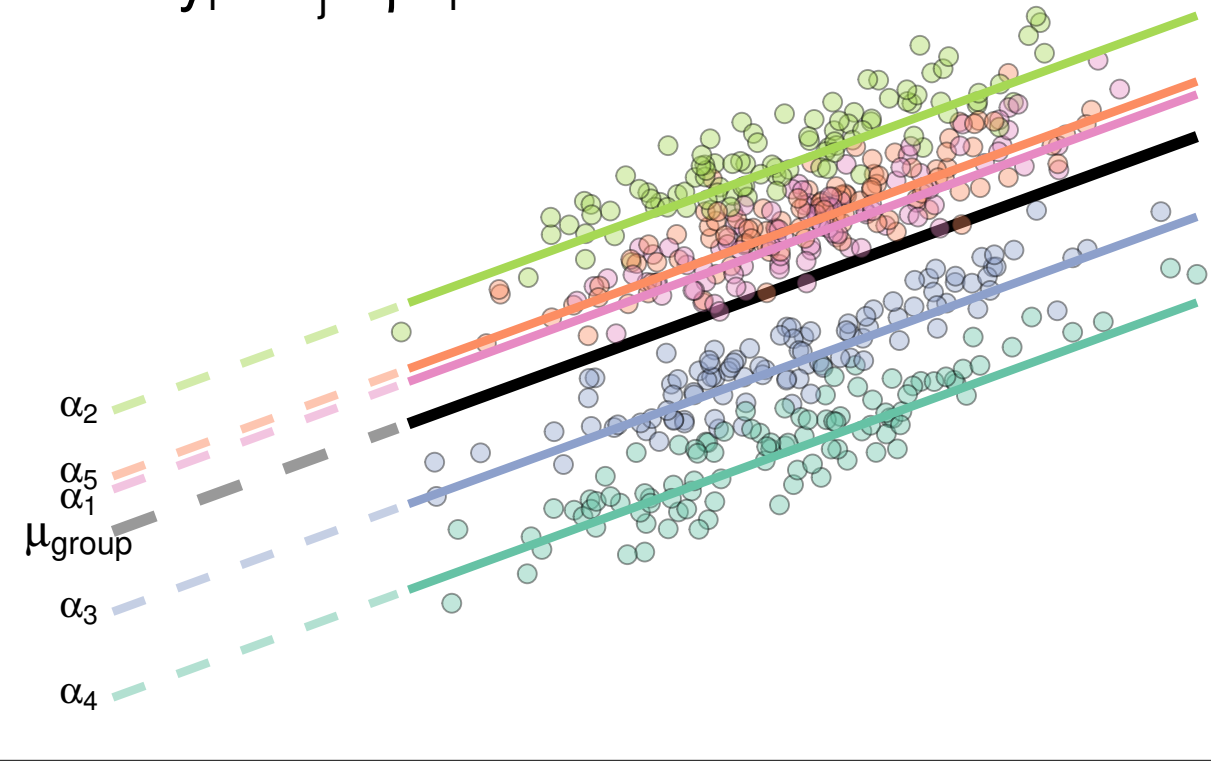Differences between Random Intercept vs Random Slope Models

(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. Point colour corresponds to group ID of the data point. The black line represents the global mean value of the distribution of random effects. (B) A random intercepts and random slopes model, where both intercepts and slopes are permitted to vary by group. Random slope models give the model far more flexibility to fit the data, but require a lot more data to obtain accurate estimates of separate slopes for each group.

**A** Random Intercepts

$$y_i = \alpha_j + \beta x_i$$

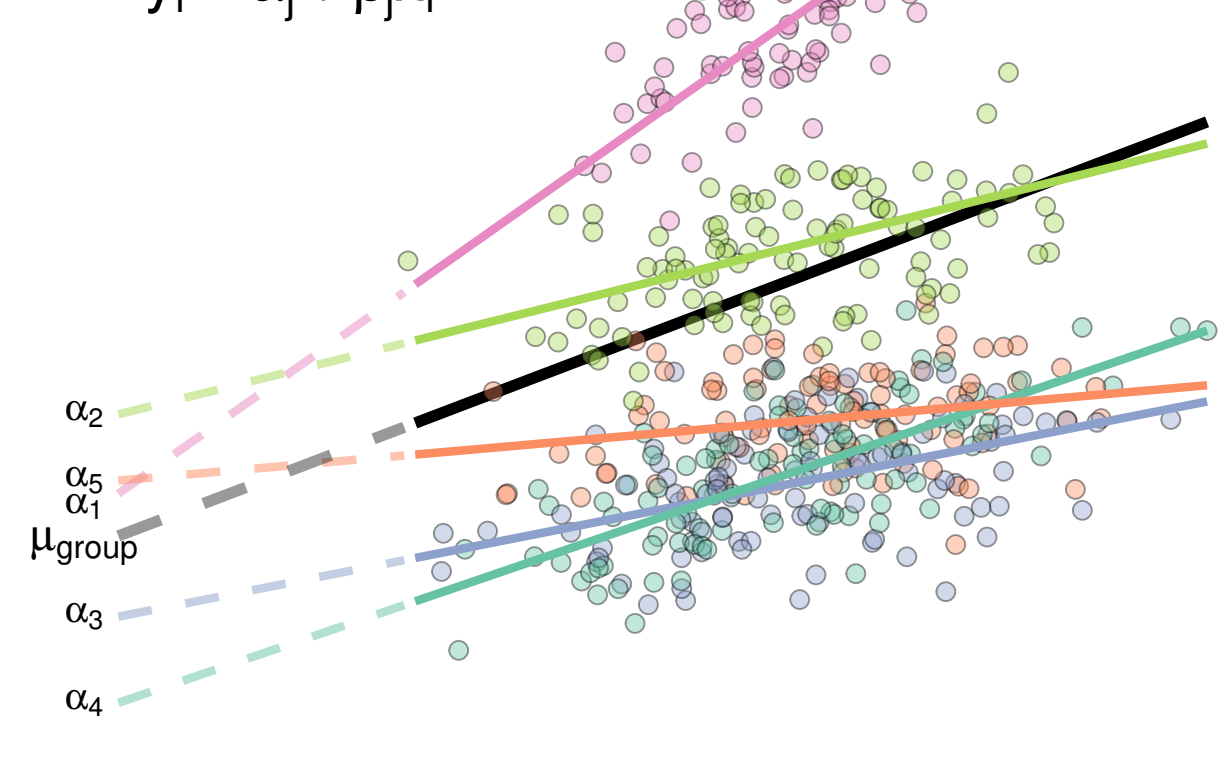**B** Random Intercepts and Slopes

$$y_i = \alpha_j + \beta_j x_i$$

**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, A) or strong correlation (r = 0.9, B). With moderate collinearity, estimation of $\beta_{x1}$ is precise, but certainty of the sign of $\beta_{x2}$ is low. When collinearity is strong, estimation of $\beta_{x1}$ is far less precise, with 14% of simulations estimating a negative coefficient for the effect of $x1$. For more elaborate versions of these simulations, see Freckleton (2011)
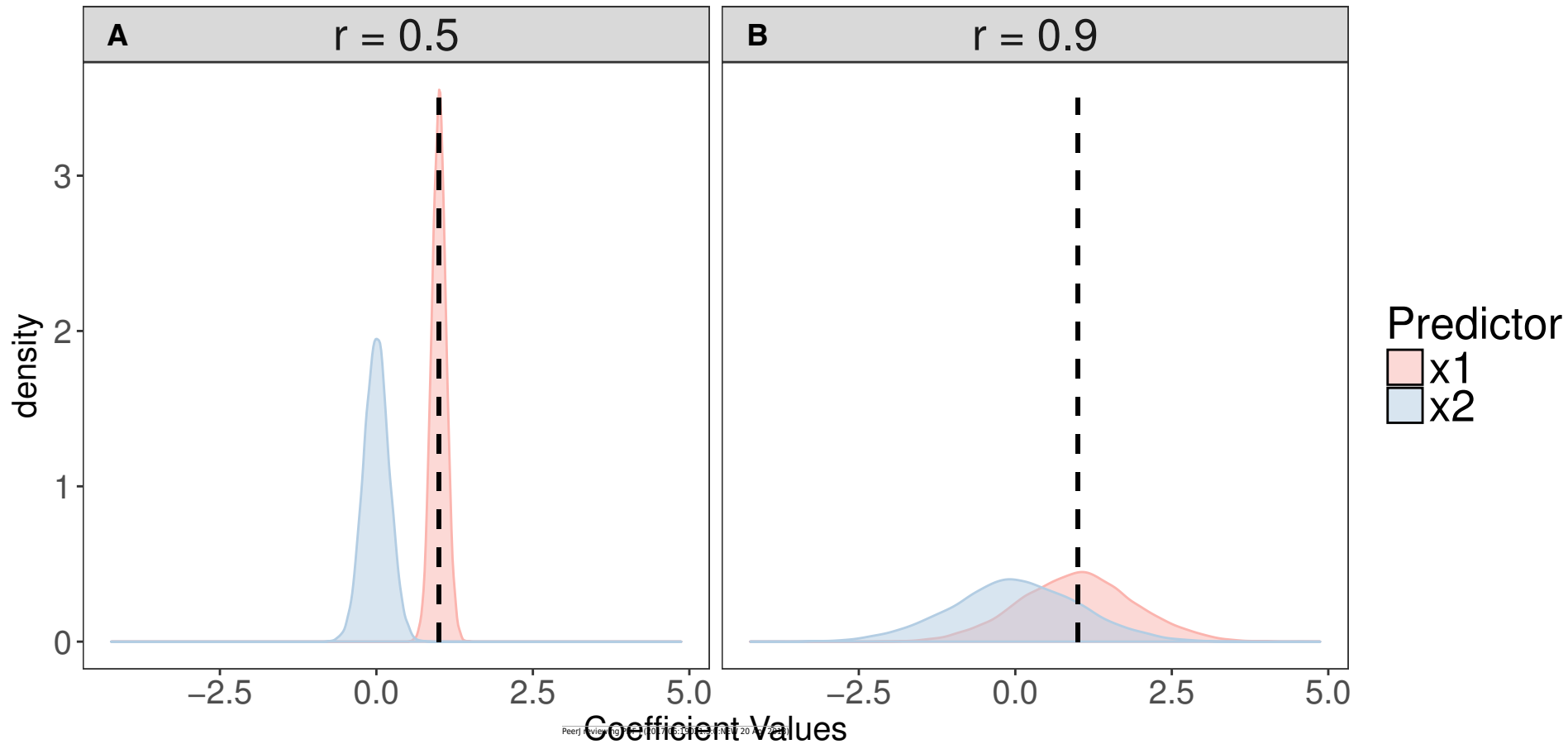
# Figure 3 (on next page)

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.