# Violating the normality assumption may be the lesser of two evils

Ulrich Knief[1,*] & Wolfgang Forstmeier[2]


[1] Division of Evolutionary Biology, Faculty of Biology, Ludwig Maximilian University of Munich, 82152 Planegg-Martinsried, Germany

[2] Department of Behavioural Ecology and Evolutionary Genetics, Max Planck Institute for Ornithology, 82319 Seewiesen, Germany

[*] Address for correspondence: Ulrich Knief, Division of Evolutionary Biology, Faculty of Biology, Ludwig Maximilian University of Munich, Grosshaderner Str. 2, 82152 Planegg-Martinsried, Germany, Phone: 0049-89-2180-74101, Fax: 0049-89-2180-74104, E-mail: knief@biologie.uni-muenchen.de

1    **Abstract**

2    **1.** Researchers are often uncertain about the extent to which it may be acceptable to violate the

3    assumption of normality of errors, which underlies the most-frequently used tests for statistical

4    significance (regression, *t*-test, ANOVA, and linear mixed models with Gaussian error).

5    **2.** Here we use Monte Carlo simulations to show that such Gaussian models are remarkably robust to

6    even the most dramatic deviations from normality.

7    **3.** We find that *P*-values are generally reliable if either the dependent variable *Y* or the predictor *X* are

8    normally distributed and that bias only occurs if both are heavily skewed (resulting in outliers in both

9    *X* and *Y*). In the latter case, judgement of significance at an α-level of 0.05 is still safe unless sample

10   size is very small. Yet, with more stringent significance criteria as is used when conducting numerous

11   tests (e.g. α = 0.0001) there is a greater risk of making erroneous judgements.

12   **4.** Generally we conclude that violating the normality assumption appears to be the lesser of two evils,

13   when compared to alternative solutions that are either unable to account for levels of non-

14   independence in the data (most non-parametric tests) or much less robust (e.g. Poisson models which

15   require control of overdispersion and sophisticated resampling). We argue that the latter may pose a

16   more substantial threat to the reliability of research findings when pragmatically acknowledging that,

17   in the majority of publications, statistical expertise is limited.

18

19   **Introduction**

20   In the biological, medical and social sciences, the validity of research findings is generally assessed

21   via statistical significance tests. Valid significance tests ensure the trustworthiness of scientific results

22   and should reduce the amount of random noise entering the scientific literature. Brunner and Austin

23   (2009) even regard this as the "primary function of statistical hypothesis testing in the discourse of

24   science". A *P*-value of $< 0.05$ is usually accepted as sufficiently low for rejecting the null hypothesis.

25   However, the validity of parametric significance tests depends on the whether model assumptions are

26   violated.

27

28    In a growing body of literature, researchers express their concerns about irreproducible results (Open

29    Science Collaboration 2015; Ebersole *et al.* 2016; Camerer *et al.* 2018; Silberzahn *et al.* 2018) and it

30    has been argued that the inappropriate use of statistics is a leading cause of irreproducible results

31    (Forstmeier, Wagenmakers & Parker 2017). Yet researchers may often be uncertain about which

32    statistical practices can be considered as safe and which are prone to yield overconfident conclusions.

33    Searching the literature, we found relatively little pragmatic advice (Box & Watson 1962; Mardia

34    1971; Lumley *et al.* 2002; Gelman & Hill 2007; O'Hara 2009; Zuur, Ieno & Elphick 2010) on the

35    question of whether and when it may be safe to violate the assumption of normality of errors, which

36    underlies the most commonly used tests for statistical significance (linear models "lm" and linear

37    mixed models "lmm" with Gaussian error, which includes the often more widely known techniques of

38    regression, *t*-test, and ANOVA). How much deviation is tolerable under which circumstances (in

39    terms of sample size and α-threshold)?

40

41    We here use Monte Carlo simulations to explore how violations of the normality assumption affect

42    the probability of drawing false-positive conclusions (the rate of type I errors), because these are the

43    greatest concern in the current reliability crisis (Open Science Collaboration 2015). We aim at

44    deriving simple rules of thumb, which researchers can use to judge whether the violation may be

45    tolerable and whether the *P*-value can be trusted. Furthermore, we provide an R package

46    ("TrustGauss") that researchers can use to explore the effect of specific distributions on the reliability

47    of *P*-values. Counter to intuition, we find that violations are rarely problematic, and we argue that the

48    commonly recommended solutions to the problem (e.g. using non-parametric tests, generalized linear

49    models) may represent a greater threat to the reliability of conclusions because of their lower

50    flexibility or robustness.

51

52    **The linear regression model and its assumptions**

53    At this point we need to briefly introduce the notation for the model of least squares linear regression.

54    In its simplest form, it can be formulated as $Y_i = a + b \times X_i + e_i$, where each element of the dependent

55    variable $Y_i$ is linearly related to the predictor $X_i$ through the regression coefficient $b$ (slope) and the

3

56    intercept $a$. $e_i$ is the error or residual term, which describes the deviations of the actual from the true

57    unobserved (error) or the predicted (residual) $Y_i$ and whose sum equals zero (Sokal & Rohlf 1995;

58    Gelman & Hill 2007). An $F$-test is usually employed for testing the significance of regression models

59    (Ali & Sharma 1996).

60

61    Basic statistics texts introduce (about) five assumptions that need to be met for interpreting all

62    estimates from linear regression models safely (validity, independence, linearity, homoscedasticity

63    and normality; Gelman & Hill 2007). Recall that these criteria are concerned with the dependent

64    variable $Y$, or — to be more precise — the regression error $e$. The predictor $X$ is usually not

65    considered. We refrain from revisiting all criteria in detail, but want to specifically focus on the

66    normality assumption here, which is usually tested via inspecting the distribution of the dependent

67    variable or of the residuals (Zuur, Ieno & Elphick 2010). Both visual approaches (probability or QQ-

68    plots) and formal statistical tests (Shapiro-Wilk) are commonly applied. Formal tests for normality

69    have been criticized because they have low power at small sample sizes and almost always yield

70    significant deviations from normality at large sample sizes (Ghasemi & Zahediasl 2012). Thus,

71    researchers are left with their intuition to decide how severely the normality assumption is violated

72    and how robust regression is to such violations.

73

74    Normally distributed errors are generally assumed to be the least important (yet probably the most

75    widely known) out of the five regression assumptions (Gelman & Hill 2007). Deviations from

76    normality do not bias regression coefficients (Williams, Grajales & Kurkiewicz 2013) and usually do

77    not impair hypothesis testing (no inflated type I error rate, e.g. Bishara & Hittner 2012; Puth,

78    Neuhauser & Ruxton 2014; Ives 2015; Szöcs & Schäfer 2015; Warton *et al.* 2016) even at relatively

79    small sample sizes, and with large sample sizes $\geq 500$ the Central Limit Theorem guarantees that the

80    test statistic is on average normally distributed (Lumley *et al.* 2002). Importantly, the robustness of

81    regression methods to deviations from normality of the regression errors $e$ does not only depend on

82    sample size, but also on the distribution of the predictor $X$ (Box & Watson 1962; Mardia 1971).

83

84 **Simulations to assess effects on *P*-values**

85 To illustrate the consequences of violating the normality assumption, we performed Monte Carlo

86 simulations on five continuous and five discrete datasets that were severely skewed, platy- and

87 leptokurtic or zero-inflated (distributions D0–D9; **Figure 1A** left column), going beyond previous

88 studies that examined less dramatic violations (Bishara & Hittner 2012; Puth, Neuhauser & Ruxton

89 2014; Ives 2015; Szöcs & Schäfer 2015; Warton *et al.* 2016). We explored these 10 distributions

90 across a range of sample sizes ($N$ = 10, 25, 50, 100, 250, 500, 1000). Starting with the normal

91 distribution D0 for reference, we sorted the remaining distributions D1–D9 by increasing tendency to

92 produce strong outliers (calculated as the average distance of the maximum or minimum from the

93 sample mean relative to the standard deviation of the sample for the case of $N$ = 100). We used these

94 data both as our dependent variable $Y$ and as our predictor variable $X$ in linear regression models,

95 yielding $10 \times 10 = 100$ combinations of $Y$ and $X$ for each sample size (see **Figure S1** for distributions

96 of the independent variable $Y$, the predictor $X$, and residuals).

97

98 We assessed the significance of the models via an $F$-test wherever possible and used a likelihood ratio

99 test otherwise. We fitted these models to 50,000 datasets for each combination of the dependent and

100 predictor variable. We did not simulate any effect, which means that both the regression coefficient $b$

101 and the intercept $a$ were on average zero. This enabled us to use the frequency of all models that

102 yielded a $P$-value $\leq 0.05$ as an estimate of the type I error rate at an α-level of 0.05. The null

103 distribution of $P$-values is uniform on the interval [0,1] and because all $P$-values are independent and

104 identically distributed, we constructed confidence intervals using a beta-distribution (cf. Casella &

105 Berger 2002; QQ-plots of expected vs observed $P$-values are depicted in **Figure S1**). We assessed the

106 deviation of observed from expected $-\log_{10}(P\text{-value})$ at an expected value of 3 ($P = 10^{-3}$) and 4 ($P =$

107 $10^{-4}$) and by estimating the scale shift parameter $\upsilon = \sigma_{observed} / \sigma_{expected}$ (Lin 1989), where σ is the

108 variance in $-\log_{10}(P\text{-value})$.

109

110 Since some of the predictor variables were binary rather than continuous, our regression models also

111 comprise the situation of classical two-sample $t$-tests, and we assume that the results would also

5

112    generalize to the situation of multiple predictor levels (ANOVA), which can be decomposed to

113    multiple binary predictors. To demonstrate that our conclusions from univariate models (involving a

114    single predictor) generalize to the multivariate case (involving several predictors), we fitted the above

115    models with a sample size of $N = 100$ to the same 10 dependent variables with three normally

116    distributed predictors and one additional predictor sampled from the 10 different distributions. We

117    further fitted the above models as mixed-effects models using the lme4 R package (v 1.1-14; Bates *et*

118    *al.* 2015). For that we simulated $N = 100$ independent samples each of which was sampled twice, such

119    that the single random effect "sample ID" explained roughly 30% of the variation in $Y$. We encourage

120    readers to try their own simulations using our R package "TrustGauss".

121

122    **Results**

123    The rate at which linear regression models with Gaussian error structure produced false-positive

124    results (type I errors) was very close to the expected value of 0.05 (**Figure 1B).** When sample size

125    was high ($N = 1000$), type I error rates ranged only between 0.044 and 0.052, across the 100

126    combinations of distributions of the dependent variable $Y$ and the predictor $X$. Hence, despite of even

127    the most dramatic violations of the normality assumption (see e.g. distributions D4 and D9 in **Figure**

128    **1A**), there was no increased risk of obtaining false-positive results. At $N = 100$, the range was still

129    remarkably narrow (0.037–0.058), and only for very low sample sizes ($N = 10$) we observed 4 out of

130    100 combinations which yielded notably elevated type I error rates in the range of 0.086 to 0.11.

131    These four cases all involved combinations of the distributions D4 and D9, which yield extreme

132    outliers. For this low sample size of $N = 10$, there were also cases where type I error rates were clearly

133    too low (down to 0.015, involving distributions D1–D3 where extreme values are rarer than under the

134    normal distribution D0; for details see **Table S1**).

135

136    Next we examine the scale shift parameter (**Figure 1C**) which evaluates the match between observed

137    and expected $P$-values across the entire range of $P$-values (not only the fraction at the 5% cut-off).

138    Whenever either the dependent variable $Y$ or the predictor $X$ was normally distributed, the observed

139    and expected $P$-values corresponded very well (first row and first column in **Figure 1C**).

140    Accordingly, the *P*-values fell within the 95% confidence bands across their entire range (rightmost

141    column in **Figures S1**). This observation was unaffected by sample size (**Table S2**). However, if both

142    the dependent variable *Y* and the predictor *X* were heavily skewed, consistently inflated *P*-values

143    outside the confidence bands occurred, yet this was almost exclusively limited to the case of $N = 10$

144    (**Figure 1C**). For larger sample sizes only the most extreme distribution D9 produced somewhat

145    unreliable *P*-values (**Figure 1C**). This latter effect of unreliable (mostly anti-conservative) *P*-values

146    was most pronounced when judgements were made at a very strict α-level (**Figure 1D** $α = 0.001$ and

147    **Figure 1E** $α = 0.0001$). At a sample size of $N = 100$, and for $α = 0.001$, observed $-\log_{10}(P\text{-values})$

148    were biased maximally 3.36-fold when both *X* and *Y* were sampled from distribution D9. This means

149    that *P*-values of about $P = 10^{-10}$ occurred at a rate of 0.001 ($P = 10^{(-3 \times 3.36)} = 10^{-10.08}$; **Figure 1D**). At *N*

150    $= 100$, and for $α = 0.0001$, the bias was maximally 4.54-fold (**Figure 1E**). Our multivariate and

151    mixed-model simulations confirmed that these patterns are general and also apply to models with

152    multiple predictor variables (**Figure S3**) and to models with random effects (**Figures S4**).

153

154    In summary, *P*-values from such Gaussian models are highly robust to even substantial violation of

155    the normality assumption and can be trusted, except when involving distributions with extreme

156    outliers (distribution D9). For very small sample sizes, judgements should preferably be made at α =

157    0.05 (rather than at more strict thresholds) and should also beware of outliers in both *X* and *Y*.

158

159    **Drawbacks of alternative solutions**

160    When the assumption of normality of errors is not met, it is often recommended to switch to either

161    non-parametric tests (e.g. Spearman rank correlation, Wilcoxon signed-rank test) or to model a more

162    specific error structure in a generalized linear model "glm" (e.g. binomial, negative binomial, Poisson,

163    zero-inflated Poisson). How risky are these approaches in terms of yielding type I errors?

164

165    In contrast to Gaussian models, for instance Poisson models are not at all robust to violations of the

166    distribution assumption. For comparison, we fitted the above univariate models involving the five

167    discrete distributions (D1, D2, D4, D6, D8) with a sample size of $N = 100$ using a Poisson error

7

168    structure. This yielded heavily biased type I error rates (at $\alpha = 0.05$) in either direction ranging from 0

169    to as high as 0.56, (**Figures S2**). Such inflations of type I error rates in glms have been reported

170    frequently (Warton & Hui 2011; Ives 2015; Szöcs & Schäfer 2015; Warton *et al.* 2016) and this

171    problem threatens the reliability of research whenever such models are implemented with insufficient

172    statistical expertise. First, it is absolutely essential to control for overdispersion in the data, which may

173    be particularly strong when Poisson errors are applied to measurements of areas (e.g. counts of pixels

174    or $mm^2$), latencies (e.g. counts of seconds), or concentrations (e.g. counts of molecules), besides the

175    more classical abundances (e.g. counts of animals). Failure to account for overdispersion will

176    typically result in very high rates of type I errors (Warton & Hui 2011; Ives 2015; Szöcs & Schäfer

177    2015; Warton *et al.* 2016; Forstmeier, Wagenmakers & Parker 2017). Second, even after accounting

178    for overdispersion, some models may still yield inflated type I error rates, therefore requiring

179    statistical testing via a resampling procedure (Warton & Hui 2011; Ives 2015; Szöcs & Schäfer 2015;

180    Warton *et al.* 2016). While most statistical experts might advocate for such a sophisticated approach

181    to count data, we are concerned about practicability when non-experts have to make decisions about

182    the most adequate resampling procedure. In this field of still developing statistical approaches it

183    seems much easier to get things wrong (and obtain a highly overconfident *P*-value) than to get

184    everything right. Finally, with the inclusion of random effects glmms are much more computationally

185    intensive than lmms and often fail to converge, leading to the recommendation to model all traits as

186    Gaussian (e.g. Ives & Garland 2014).

187

188    The biggest downside of non-parametric approaches is that they are less advanced and user-friendly

189    compared to linear (mixed) models (e.g. Akritas & Brunner 2003), such that only simple procedures

190    are widely known and applied. The latter, however, are applicable only to the simplest and idealized

191    scenario of fully independent data points and of only a single explanatory variable with no

192    confounding factors or covariates to be controlled for. Real data sets rarely fulfil that condition, such

193    that simple non-parametric tests often suffer from pseudoreplication and unaccounted confounds.

194    Pseudoreplication, i.e. overestimation of the number of truly independent replicates, results in

195    overconfident estimates and hence is one of the leading causes of false-positive conclusions

196     (Forstmeier, Wagenmakers & Parker 2017). Gaussian models, in contrast, allow us to easily control

197     for pseudoreplication by specifying the random effects that cause non-independence of data points

198     (mixed-effects models).

199

200     Finally, there is much to be gained when researchers follow a standardized way of reporting effect

201     sizes (Lumley *et al.* 2002). For instance, a study that examines the effect of a single treatment on

202     multiple dependent variables (e.g. health parameters) may often switch forth and back between

203     reporting parametric and non-parametric test statistics depending on how strongly the trait of interest

204     deviates from normality, rendering a comparison of effect sizes difficult. Methods of converting effect

205     sizes for discrete traits (e.g. odds ratio from a 2×2 contingency table) into effect sizes for continuous

206     traits (e.g. Pearson correlation coefficient) already work by violating the normality assumption (e.g.

207     fitting a Pearson correlation through the binary data of a 2×2 table; Nakagawa & Cuthill 2007), so

208     why not always report the Gaussian model to begin with, if the primary purpose of the test is to obtain

209     a reliable *P*-value?

210

211     **Practical advice (for referees)**

212     In order to effectively guard against false-positive claims entering the scientific literature, violations

213     of the normality assumption in linear models are much less of a problem than violations of the

214     independence of data points (pseudoreplication; Schielzeth & Forstmeier 2009; Forstmeier,

215     Wagenmakers & Parker 2017). To avoid the negative consequences of strong deviations from

216     normality that may occur under some conditions (see **Figure 1**) it may be most advisable to apply a

217     rank-based inverse normal (RIN) transformation (aka rankit scores Bliss 1967) to the data, which can

218     approximately normalize most distributional shapes and which effectively minimizes type I errors and

219     maximises statistical power (Bishara & Hittner 2012).

220

221     In practice, we recommend the following to referees:

222     (1) When a test assumes Gaussian errors, request a check for outliers, particularly if very small *P*-

223     values are reported. Consider recommending a RIN-transformation for strong deviations from

224     normality.

225     (2) For Poisson, binomial and negative binomial errors, always check whether the issues of

226     overdispersion and resampling are addressed, otherwise request an adequate control for type I errors

227     or verification with Gaussian models.

228     (3) Requesting a switch to non-parametric statistics is not advised, and requests for switching from lm

229     to glm (or from lmm to glmm) should be accompanied with sufficient advice (e.g. R-code) to ensure a

230     safe implementation.

231

232     **Conclusion**

233     If we are interested in statistical hypothesis testing, linear regression models with Gaussian error

234     structure are generally robust to violations of the normality assumption. Judging *P*-values at the

235     threshold of $\alpha = 0.05$ is nearly always safe, but if both *Y* and *X* are skewed, we should avoid being

236     overly confident in very small *P*-values and examine whether these result from outliers in both *X* and

237     *Y* (see also Osborne & Overbay 2004). With this caveat in mind, violating the normality assumption is

238     relatively unproblematic. Alternative solutions like Poisson models and non-parametric tests may bear

239     a greater risk of yielding anti-conservative *P*-values when applied by scientists with limited statistical

240     expertise.

241

242     **Data availability**

243     All functions are bundled in an R package named "TrustGauss". The R package and all

244     supplementary figures are accessible through the Open Science Framework (osf.io/r5ym4).

245

246     **Acknowledgements**

249

250 **Author contributions**

251 WF and UK conceived of the study. UK wrote the simulation code. UK and WF prepared the

252 manuscript.

253

254 **Competing interests**
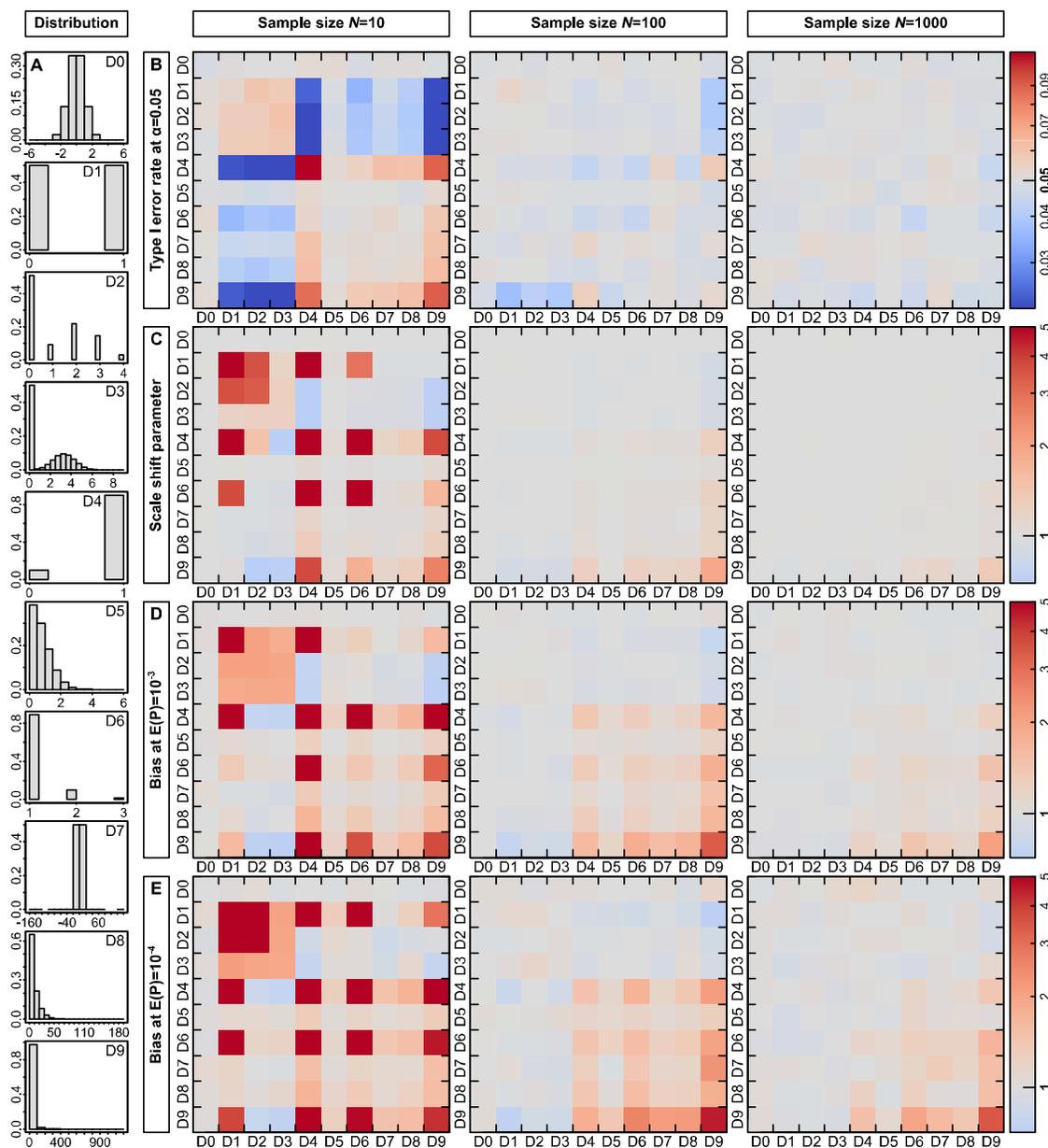
255 The authors declare no competing financial interests.

256

257 **References**

258 Akritas, M.G. & Brunner, E. (2003) Nonparametric models for ANOVA and ANCOVA: a review.
259     *Recent Advances and Trends in Nonparametric Statistics*, 79–91.
260 Ali, M.M. & Sharma, S.C. (1996) Robustness to nonnormality of regression F-tests. *Journal of*
261     *Econometrics, 71,* 175–205.
262 Bates, D., Mächler, M., Bolker, B.M. & Walker, S.C. (2015) Fitting linear mixed-effects models
263     using lme4. *Journal of Statistical Software, 67,* 1–48.
264 Bishara, A.J. & Hittner, J.B. (2012) Testing the significance of a correlation with nonnormal data:
265     comparison of Pearson, Spearman, transformation, and resampling approaches. *Psychological*
266     *Methods, 17,* 399–417.
267 Bliss, C.I. (1967) *Statistics in biology*. McGraw-Hill, New York, NY.
268 Box, G.E.P. & Watson, G.S. (1962) Robustness to non-normality of regression tests. *Biometrika, 49,*
269     93–106.
270 Brunner, J. & Austin, P.C. (2009) Inflation of type I error rate in multiple regression when
271     independent variables are measured with error. *Canadian Journal of Statistics, 37,* 33–46.
272 Camerer, C.F., Dreber, A., Holzmeister, F., Ho, T.H., Huber, J., Johannesson, M., Kirchler, M., Nave,
273     G., Nosek, B.A., Pfeiffer, T., Altmejd, A., Buttrick, N., Chan, T.Z., Chen, Y.L., Forsell, E.,
274     Gampa, A., Heikensten, E., Hummer, L., Imai, T., Isaksson, S., Manfredi, D., Rose, J.,
275     Wagenmakers, E.J. & Wu, H. (2018) Evaluating the replicability of social science
276     experiments in *Nature* and *Science* between 2010 and 2015. *Nature Human Behaviour, 2,*
277     637–644.
278 Casella, G. & Berger, R.L. (2002) *Statistical inference,* 2nd edn. Duxbury Press, Pacific Grove,
279     California.
280 Ebersole, C.R., Atherton, O.E., Belanger, A.L., Skulborstad, H.M., Allen, J.M., Banks, J.B., Baranski,
281     E., Bernstein, M.J., Bonfiglio, D.B.V., Boucher, L., Brown, E.R., Budiman, N.I., Cairo, A.H.,
282     Capaldi, C.A., Chartier, C.R., Chung, J.M., Cicero, D.C., Coleman, J.A., Conway, J.G.,
283     Davis, W.E., Devos, T., Fletcher, M.M., German, K., Grahe, J.E., Hermann, A.D., Hicks,
284     J.A., Honeycutt, N., Humphrey, B., Janus, M., Johnson, D.J., Joy-Gaba, J.A., Juzeler, H.,
285     Keres, A., Kinney, D., Kirshenbaum, J., Klein, R.A., Lucas, R.E., Lustgraaf, C.J.N., Martin,
286     D., Menon, M., Metzger, M., Moloney, J.M., Morse, P.J., Prislin, R., Razza, T., Re, D.E.,
287     Rule, N.O., Sacco, T.F., Sauerberger, K., Shrider, E., Shultz, M., Siemsen, C., Sobocko, K.,
288     Sternglanz, R.W., Summerville, A., Tskhay, K.O., van Allen, Z., Vaughn, L.A., Walker, R.J.,
289     Weinberg, A., Wilson, J.P., Wirth, J.H., Wortman, J. & Nosek, B.A. (2016) Many labs 3:
290     evaluating participant pool quality across the academic semester via replication. *Journal of*
291     *Experimental Social Psychology, 67,* 68–82.
292 Forstmeier, W., Wagenmakers, E.J. & Parker, T.H. (2017) Detecting and avoiding likely false-
293     positive findings – a practical guide. *Biological Reviews, 92,* 1941–1968.
294 Gelman, A. & Hill, J. (2007) *Data analysis using regression and multilevel/hierarchical models,* 1
295     edn. Cambridge University Press, New York.

11

296 Ghasemi, A. & Zahediasl, S. (2012) Normality tests for statistical analysis: a guide for non-
297     statisticians. *Int J Endocrinol Metab,* **10,** 486–489.
298 Ives, A.R. (2015) For testing the significance of regression coefficients, go ahead and log-transform
299     count data. *Methods in Ecology and Evolution,* **6,** 828–835.
300 Ives, A.R. & Garland, T. (2014) Phylogenetic regression for binary dependent variables. *Modern*
301     *phylogenetic comparative methods and their application in evolutionary biology* (ed. L.Z.
302     Garamszegi), pp. 231–261. Springer, Berlin, Heidelberg.
303 Lin, L.I. (1989) A concordance correlation-coefficient to evaluate reproducibility. *Biometrics,* **45,**
304     255–268.
305 Lumley, T., Diehr, P., Emerson, S. & Chen, L. (2002) The importance of the normality assumption in
306     large public health data sets. *Annual Review of Public Health,* **23,** 151–169.
307 Mardia, K.V. (1971) The effect of nonnormality on some multivariate tests and robustness to
308     nonnormality in the linear model. *Biometrika,* **58,** 105–121.
309 Nakagawa, S. & Cuthill, I.C. (2007) Effect size, confidence interval and statistical significance: a
310     practical guide for biologists. *Biological Reviews,* **82,** 591–605.
311 O'Hara, R.B. (2009) How to make models add up—a primer on GLMMs. *Annales Zoologici Fennici,*
312     **46,** 124–137.
313 Open Science Collaboration (2015) Estimating the reproducibility of psychological science. *Science,*
314     **349,** aac4716.
315 Osborne, J.W. & Overbay, A. (2004) The power of outliers (and why researchers should ALWAYS
316     check for them). *Practical Assessment, Research & Evaluation,* **9,** Available online:
317     http://PAREonline.net/getvn.asp?v=9&n=6.
318 Puth, M.T., Neuhauser, M. & Ruxton, G.D. (2014) Effective use of Pearson's product-moment
319     correlation coefficient. *Animal Behaviour,* **93,** 183–189.
320 Schielzeth, H. & Forstmeier, W. (2009) Conclusions beyond support: overconfident estimates in
321     mixed models. *Behavioral Ecology,* **20,** 416–420.
322 Silberzahn, R., Uhlmann, E.L., Martin, D.P., Anselmi, P., Aust, F., Awtrey, E., Bahník, Š., Bai, F.,
323     Bannard, C., Bonnier, E., Carlsson, R., Cheung, F., Christensen, G., Clay, R., Craig, M.A.,
324     Dalla Rosa, A., Dam, L., Evans, M.H., Flores Cervantes, I., Fong, N., Gamez-Djokic, M.,
325     Glenz, A., Gordon-McKeon, S., Heaton, T.J., Hederos, K., Heene, M., Hofelich Mohr, A.J.,
326     Högden, F., Hui, K., Johannesson, M., Kalodimos, J., Kaszubowski, E., Kennedy, D.M., Lei,
327     R., Lindsay, T.A., Liverani, S., Madan, C.R., Molden, D., Molleman, E., Morey, R.D.,
328     Mulder, L.B., Nijstad, B.R., Pope, N.G., Pope, B., Prenoveau, J.M., Rink, F., Robusto, E.,
329     Roderique, H., Sandberg, A., Schlüter, E., Schönbrodt, F.D., Sherman, M.F., Sommer, S.A.,
330     Sotak, K., Spain, S., Spörlein, C., Stafford, T., Stefanutti, L., Tauber, S., Ullrich, J., Vianello,
331     M., Wagenmakers, E.-J., Witkowiak, M., Yoon, S. & Nosek, B.A. (2018) Many analysts, one
332     data set: making transparent how variations in analytic choices affect results. *Advances in*
333     *Methods and Practices in Psychological Science,* **1,** 337–356.
334 Sokal, R.R. & Rohlf, F.J. (1995) *Biometry*. W. H. Freeman, New York.
335 Szöcs, E. & Schäfer, R.B. (2015) Ecotoxicology is not normal. *Environmental Science and Pollution*
336     *Research,* **22,** 13990–13999.
337 Warton, D.I. & Hui, F.K.C. (2011) The arcsine is asinine: the analysis of proportions in ecology.
338     *Ecology,* **92,** 3–10.
339 Warton, D.I., Lyons, M., Stoklosa, J. & Ives, A.R. (2016) Three points to consider when choosing a
340     LM or GLM test for count data. *Methods in Ecology and Evolution,* **7,** 882–890.
341 Williams, M.N., Grajales, C.A.G. & Kurkiewicz, D. (2013) Assumptions of multiple regression:
342     correcting two misconceptions. *Practical Assessment, Research & Evaluation,* **18,** Available
343     online: http://pareonline.net/getvn.asp?v=18&n=11.
344 Zuur, A.F., Ieno, E.N. & Elphick, C.S. (2010) A protocol for data exploration to avoid common
345     statistical problems. *Methods in Ecology and Evolution,* **1,** 3–14.

346

347
348 **Figure 1 |** *P*-values from Gaussian linear regression models are in most cases unbiased. (**A**) Overview

349 of the ten different distributions that we simulated. Distributions D0 is Gaussian and all remaining

350 distributions are sorted by their tendency to produce strong outliers. Distributions D1, D2, D4, D6 and

351 D8 are discrete. The roman numbers refer to the plots in (**B–E**) where on the *Y*-axis the distribution of

352 the dependent variable and on the *X*-axis of the predictor is indicated. (**B**) Type I error rate at an α-

353 level of 0.05 for sample sizes of *N* = 10, 100 and 1000. Red colours represent increased and blue

354 conservative type I error rates. (**C**) Scale shift parameter, (**D**) the bias in *P*-values at an expected *P*-

355 value of $10^{-3}$ and (**E**) the bias in *P*-values at an expected *P*-value of $10^{-4}$.

13