



UNIVERSITY OF
PLYMOUTH



Peninsula Medical School
Faculty of Health

2017-12-01

A comparison of inferential analysis methods for multilevel studies: Implications for drawing conclusions in animal welfare science

KN Stevens

L Asher

K Griffin

M Friel

N O'Connell

Let us know how access to this document benefits you

General rights

All content in PEARL is protected by copyright law. Author manuscripts are made available in accordance with publisher policies. Please cite only the published version using the details provided on the item record or document. In the absence of an open licence (e.g. Creative Commons), permissions for further reuse of content should be sought from the publisher or author.

Take down policy

If you believe that this document breaches copyright please [contact the library](#) providing details, and we will remove access to the work immediately and investigate your claim.

Follow this and additional works at: <https://pearl.plymouth.ac.uk/pms-research>

Recommended Citation

Stevens, K., Asher, L., Griffin, K., Friel, M., & O'Connell, N. (2017) 'A comparison of inferential analysis methods for multilevel studies: Implications for drawing conclusions in animal welfare science', *Applied Animal Behaviour Science*, 197, pp. 101-111. Elsevier: Available at: <https://doi.org/10.1016/j.applanim.2017.08.002>

This Article is brought to you for free and open access by the Faculty of Health at PEARL. It has been accepted for inclusion in Peninsula Medical School by an authorized administrator of PEARL. For more information, please contact openresearch@plymouth.ac.uk.

2017-12

A comparison of inferential analysis methods for multilevel studies: Implications for drawing conclusions in animal welfare science

Stevens, KN

<http://hdl.handle.net/10026.1/9900>

10.1016/j.applanim.2017.08.002

Applied Animal Behaviour Science

Elsevier

All content in PEARL is protected by copyright law. Author manuscripts are made available in accordance with publisher policies. Please cite only the published version using the details provided on the item record or document. In the absence of an open licence (e.g. Creative Commons), permissions for further reuse of content should be sought from the publisher or author.

20 **ABSTRACT**

21 Investigations comparing the behaviour and welfare of animals in different environments have led to
22 mixed and often conflicting results. These could arise from genuine differences in welfare, poor
23 validity of indicators, low statistical power, publication bias, or inappropriate statistical analysis. Our
24 aim was to investigate the effects of using four approaches for inferential analysis of datasets of
25 varying size on model outcomes and potential conclusions. We considered aggression in 864 growing
26 pigs over six weeks as measured by ear and body injury score and relationships with: less and more
27 enriched environments, pig's relative weight, and sex. Pigs were housed in groups of 18 in one of four
28 pens, replicating the experiment 12 times. We applied four inferential models that either used a
29 summary statistic approach, or else fully or partially accounted for complexities in study design. We
30 tested models using both the full dataset ($n = 864$) and also using small sample sizes ($n = 72$).

31 The most appropriate inferential model was a mixed effects, repeated measures model to compare ear
32 and body score. Statistical models that did not account for the correlation between repeated measures
33 and/or the random effects from replications and pens led to spurious associations between
34 environmental factors and indicators of aggression, which were not supported by the initial
35 exploratory analysis. For analyses on smaller datasets ($n = 72$), due to the effect size and number of
36 independent factors, there was insufficient power to determine statistically significant associations.

37 Based on the mixed effects, repeated measures models, higher body injury scores were associated
38 with more enrichment (coef. est. = 0.09, $p = 0.02$); weight (coef. est. = 0.05, $p < 0.001$); pen location
39 on the right side (coef. est. = 0.08, $p = 0.03$) and at the front of the experimental room (coef. est. =
40 0.11, $p = 0.003$). By comparison, lower ear injury scores were associated with more enrichment
41 (coef. est. = -0.51, $p = 0.005$) and pen location at the front of the experimental room (coef. est. = -0.4,
42 $p = 0.02$). These observed differences support the hypothesis that injuries to the body and ears arise
43 from different risk factors. Although calculation of the minimum required sample size prior to
44 conducting an experiment and selection of the inferential analysis method will contribute to the

45 validity of the study results, conflict between the outcomes will require further investigation via
46 different methods such as sensitivity and specificity analysis.

47

48 Word count: 400

49

50 **Key Words:** Study design, sample size, mixed effect models, pig, animal health, animal welfare.

51

52 1. INTRODUCTION

53 The statistician George Box stated “all models are wrong, but some are useful” (Box and Draper,
54 1987); which raises the question, how do we determine which statistical model, or in other
55 terminology, inferential analysis method, is most appropriate? In recent years, a spotlight has been
56 directed at the transparency of animal research methodology, with low rates of methodological
57 reporting being associated with less scientific rigour and lower reproducibility (Vogt et al 2016,
58 Ionnides et al 2009, Kilkenny et al 2009). Articles pertaining to animal research have been criticised
59 in the past for their design, statistical analysis and reporting (McCance, 1995; Kilkenny et al., 2009;
60 Sargeant et al., 2010). The publication of a list of guidelines for animal research known as the
61 ARRIVE guidelines (Kilkenny et al., 2010), has helped to improve the quality of animal research
62 (Gulin et al., 2015). These guidelines highlight the importance of choosing the appropriate
63 experimental assessments, sample sizes and statistical inferential analysis methods. It is important to
64 ensure the sample size is sufficient to test the study hypothesis, but also bearing in mind the ethical
65 and financial implications of using an unnecessarily large sample size within an experiment. There is
66 a plethora of techniques to produce sample size estimates, and the appropriate technique will depend
67 on the inferential analysis used for a study. Sample size can often be quite difficult to calculate for
68 more complex designs, though the importance of conducting these calculations accurately has been
69 well communicated, particularly in clinical trials literature (Freiman et al., 1978; Biau et al., 2008).

70 Discussion in this area naturally leads into consideration of the methodology of the statistical analysis
71 conducted on the collected data. Many of the papers focussing on the quality of research using
72 animals have primarily targeted experimental design, animal numbers, and reporting, but have not
73 discussed the appropriate analysis of what can often be complex datasets. Precise replication of a
74 published study is rarely performed, and typically different studies will use different experimental
75 designs and statistical inferential techniques to address the question. Although this can make
76 comparisons between published studies difficult, agreement in the overall conclusions under such
77 circumstances can be considered strong evidence for the named association, though more subtle or
78 complex relationships may potentially be missed. An identified significant treatment effect across

79 studies through use of meta-analysis, is typically considered to be robust evidence for an association,
80 and also allows the magnitude of the effect size to be more precisely estimated than in single studies
81 considered in isolation (Borenstein et al., 2009). However meta-analysis also has limitations, for
82 example when few studies have been published in an area, when they differ substantially, or when the
83 inferential analysis used is inappropriate for the design.

84 Within the field of animal welfare, many published results on a particular issue are mixed or
85 conflicting, leading to somewhat mixed messages about what the most appropriate solution for an
86 identified welfare hazard might be. To some extent, it is possible that this is at least partly due to
87 publication bias (e.g. Hopewell et al., 2009; Brown et al., 2017) and the drive for novelty rather than
88 further support for a set of hypotheses in published research. However, the lack of agreement
89 between studies may be due to other factors – the differences may reflect genuine differences between
90 the studies, arising for reasons as yet unmeasured or unaccounted for. They may be due to the use of
91 indicators that have not been thoroughly validated in all respects for the species in question (Cronbach
92 & Meehl, 1955). Finally, the observed lack of agreement may be due to inappropriate statistical
93 analysis, leading to masking of true effects, or the discovery of false positives.

94 Even when two studies ask a very similar research question with largely similar methodology, mixed
95 results can emerge. A typical example of this can be found in studies that investigate causes, and
96 consequently solutions, for aggression in pigs. For example, Beattie et al. (1996) investigated whether
97 an enrichment object or floor space had more influence on pig behaviour. Their analysis showed that
98 duration of harmful behaviour was significantly higher in less enriched pens, and measured pig
99 aggressive behaviours had no significant association with space allowance. By comparison, Turner et
100 al. (2000) found that smaller space allowances were associated with more skin lesions and longer-
101 lasting aggressive events. These studies were similar in a number of respects, except that Turner et al.
102 (2000) regularly adjusted pen sizes to maintain a consistent stocking density (weight per m²)
103 throughout the experiment, whereas Beattie et al. (1996) maintained pen dimensions (hence stocking
104 density would increase throughout the study). Consequently, the two studies are incomparable with
105 conventional meta-analytic approaches. Variation in the indicators used could also potentially explain

106 differences in model outcomes For example, different indicators of injuries in pigs result in
107 differences in the final conclusion, even if the studies use otherwise similar experimental designs and
108 methods for inferential analysis. In relation to the provision of straw for pigs, different indicators of
109 aggression have lead to different conclusions; for example, Lahrmann et al. (2015) found reduced
110 shoulder injuries for straw-housed pigs, whereas Morgan et al. (1998) found that straw-housed pigs
111 performed more aggressive interactions and Statham et al. (2011) and Arey and Franklin (1995) have
112 both reported no significant effect of the provision of straw on outbreaks of aggression. Aggression
113 can, and indeed, has been described and measured using a wide variety of indicators. Examples of
114 indicators for aggression are: duration of fights and number of bites (Andersen et al. (2000);
115 prevalence of giving/ receiving belly nosing, mounting, ear and tail biting, and biting the pen bars,
116 chains or other pen details (Brunberg et al. (2011); the ratio of aggressive events to social interactions
117 (Drickamer et al., 1999); skin lesions on different body areas (Desire et al., 2016). Frequently, there is
118 little or no overlap between studies, or construct validation to demonstrate that all indicators recorded
119 measure what they are proposed to measure (e.g. tail biting has been considered an indicator of
120 aggression; however this has been reconsidered in more recent years, e.g. Taylor et al., 2010).

121 Here we used a study investigating aggression in pigs to compare differences between two areas for
122 the assessment of skin injuries (believed to be indicative of aggression in pigs), an ear score and a
123 composite body score (Conte et al. 2012), and the effects of analysing the data via four inferential
124 methods: (i) generalised linear models; (ii) repeated measures analysis; (iii) linear mixed effect
125 models; and (iv) linear mixed effect models for repeated measures. We compare the significant
126 associations between the two injury assessments and the covariates detected via the exploratory and
127 four methods of inferential analysis. These four approaches were chosen because, to varying degrees,
128 these models could account for some of the features of the data and model parameters could be
129 directly interpreted.

130 Methods (i)-(iii) were considered sub-optimal relative to (iv), as these models were unable to account
131 for correlation in the repeated measures, and /or random effects from the hierarchical structure in the
132 data (pens within replication). We hypothesised that not accounting for random effects from the pens

133 within replication and correlation between repeated measures will either result in additional spurious
134 relationships and/or mask possible significant relationships between our injury assessments and the
135 covariates. By ignoring random effects, we hypothesise there will be more statistically significant
136 associations with environmental factors, and by ignoring the repeated measurements, we hypothesise
137 the association between injury score and time covariate will be more complex.

138 We investigated the effects of sample size within multilevel designs by analysing the data from
139 different replications (n=18 pigs * 4 pens per replicate) as separate studies, and comparing the
140 coefficient estimates from each of these analyses. A reduced sample size leads to a decrease in power,
141 which means it is more difficult to identify the environmental factors associated with the injury
142 scores. We hypothesize, that with a reduced sample size, there will be fewer statistically significant
143 associations between injury scores and environmental factors.

144

145 **2. METHODS AND MATERIALS**

146 *2.1 Animals and Housing*

147 The study was conducted at the Agri-Food and Biosciences Institute, Hillsborough, County Down,
148 Northern Ireland. The study used commercial crossbreed PIC 337 (Large White x Landrace) pigs.
149 Pigs received a commercial weaner diet ad libitum and water was always available, according to the
150 standard practices on the farm.

151 Each pig was weighed when they were four weeks and again at ten weeks old. The pigs' sex and
152 weights at 4 weeks of age were used by the stockman to balance the groups to achieve a similar
153 average weight and 50:50 sex ratio in each group of 18 individuals. Groups were then allocated at
154 random to one of four pens. The pigs remained in these pens for a period of approximately six weeks,
155 and the study was replicated twelve times, which led to a sample size of 864.

156 Pigs were assigned to one of four pens for the study that were contained within an experimental room
157 situated in a long shed, which was divided into a series of similar rooms, with floor to ceiling solid
158 walls between each room. Two types of pen environment were used within this study. Pens 1 and 3
159 were classed as more enriched environments; these pens were 2.18 m × 5.16 m in dimension with
160 deep straw bedding (replenished weekly). Pens 2 and 4 were classed as less enriched environments,
161 these were 2.18 m × 3.42 m in dimension, and no straw was provided. Both pens had floors
162 constructed from concrete and were partially slatted, however in the more enriched pens (1 and 3) the
163 slats were covered with plywood to prevent straw falling into the slurry system. In all pens, suspended
164 wooden blocks were provided as standard enrichment.

165 Pens 1 and 2 were located on the left side of the experimental room and pens 3 and 4 were located on
166 the right. The adjacent room on the right (next to pens 3 and 4) almost always contained weaner pigs,
167 whereas the adjacent room on the left (next to pens 1 and 2) was frequently empty, or was
168 occasionally used to house sows that could not enter farrowing crates. The difference in directional
169 noise from each adjacent room was balanced in the experimental design by having one pen of each
170 treatment type on both sides of the room. Two of the four pens were located next to the front of the
171 room (pen 2 and pen 3), and the other two pens were located at the back next to an internal corridor.

172 The pigs were kept commercially, hence decisions relating to culling and health were made by the
173 farm manager, as part of the standard on-farm procedures. Outbreaks of aggression leading to injury
174 were observed only on video footage, analysed typically several weeks after recording took place.
175 Animals that were observed to have high body scores were reported to farm staff, and monitored
176 closely by farm staff and researchers for a period of 7 days after. No animals were culled for the
177 purposes of this study, though as noted in section 2.3, a small number of animals (n=9 out of 862
178 pigs) died during the study period due to poor health or failure to thrive.

179 ***2.2 Assessment of Injury***

180 An assessment of each individual's injuries was completed at three time points after entering the pens:
181 (1) On day 4; (2) Between days 8 – 17; (3) Between days 29 and 39. At each assessment each pig was

182 scored on the following body areas: left and right ear; snout; left and right shoulder; front and back
183 legs; left and right flank; left and right hindquarter; and back; using a six point scaling system, as
184 defined in figure 1 (Conte et al. 2012). As part of the standard practice on the farm, 50% of the tail
185 was docked within the first 24 hours after birth for every pig, this meant that tail score had limited
186 value as an indicator for aggression.

187 *2.2.1 Indicators of Aggression*

188 Ear and body score were considered as indicators of aggression. At each assessment time point, the
189 ear score was recorded as the higher observed injury score on either the left or right ear (possible
190 score 0-5), and the body score was recorded as the sum score of the back, left and right shoulder,
191 flank and hindquarters scores (possible score 0 – 25).

192 Due to the method used to construct the body score, based on the Conte et al (2012) scale, the two
193 elements of frequency of injury and severity are confounded, especially for lower values. In our
194 dataset, body score ranged between zero and 25, suggesting body score could be analysed as a
195 continuous variable. A histogram plot of the log transformed body score implied we could assume the
196 data followed a Gaussian distribution.

197 Each ear was scored on a scale between zero and five, with a score of zero signifying no injuries or
198 damage, and a score of five indicating the presence of many deep red lesions. As very few pigs were
199 identified with a score of 3 or more, categories 3 to 5 were combined, so that the ear score categories
200 represented: 0 = no injuries; 1 = one small superficial lesion; 2 = more than one small, superficial
201 lesion; or one red (ie deeper than score 1) superficial lesion; 3 = one or more deep lesions, or more
202 than one red superficial lesions. Initial exploratory analysis suggested that the relationship between
203 the housing conditions, sex and weight were similar for pigs with an ear score of 0 or 1. Therefore,
204 these two groups were combined to simplify subsequent inferential analyses.

205

206

207 **2.3 Statistical Analysis**

208 As injury assessments were made at three irregularly spaced points in time, the assessments for an
209 individual pig could be correlated, but the strength of the correlation may differ because of the
210 variable time differences. Replicating the study 12 times may cause significant random effects for
211 each pen within replication. The differences could be caused by the combination of pigs within a pen,
212 or even associated with unmeasured external influences (e.g. weather conditions, handler behaviour,
213 noise). Using weight at 4 and 10 weeks of age, we produced estimates of each individual's
214 intermediate weights by fitting a linear model between the two time points. Although growth is
215 usually statistically modelled by a curve, plots of the expected growth curves in Carr (1998) indicated
216 that a linear estimate of pig weight would be an appropriate approximation over the short time scale
217 used in this study.

218 We calculated individual relative weights in each pen within replication, in line with previous
219 research indicating that an individual's relative size compared with its group mates is more important
220 than its actual size (Nettle et al., 2013). Andersen et al. (2000) found no significant difference in
221 number of bites between groups of pigs with low and high weight variability, which suggested
222 removing pen differences would have no adverse effects. This is similar to comparing a pig's weight
223 rank, but also accounts for variable weight differences between pigs.

224 Missing data were due to human error in data entry, and death or culling of the individual pig during
225 the course of the study, either due to poor health or failure to thrive.

226 The plots and statistical analyses were produced using the statistical program R (Team, 2015) using
227 the multgee (Touloumis, 2016), ordinal (Christensen, 2015), and lme4 (Bates et al., 2015) packages to
228 produce the statistical models.

229 **2.3.2 Exploratory Analysis**

230 Before applying any statistical test or fitting a statistical model to data, it is important to perform
231 appropriate exploratory analysis. Choosing the right method to explore the data will depend on the

232 question being addressed. As these data consisted of observations measured over time, we aimed to
233 explore how body and ear score changed over time.

234 We plotted each pig's body score over time and fitted a Gaussian kernel smooth estimator to pigs
235 within each category (i.e. by treatment enrichment level). A kernel estimator is a non-parametric
236 method of fitting a line between two continuous variables. If there is uncertainty about the form of
237 this relationship (i.e. linear, quadratic, etc.), visual inspection of plots of the data can provide insight
238 into this. An appropriate bandwidth is determined, with bigger bandwidths creating smoother lines.
239 We selected a bandwidth of 15, as injury assessments took place every 14 days on average (more
240 details of kernel estimators can be found in Wand and Jones (1994)). As we were treating ear score as
241 an ordinal variable, we looked at the proportional change of pigs within each category, and used the
242 same methods as outlined above for body score.

243 *2.3.3 Inferential Analysis*

244 The data from this experiment possessed a hierarchical structure, where we had repeated
245 measurements for each pig, within a pen, within a replication. There are various methods that can be
246 applied to this type of data, depending on the assumptions one makes. We compared the results of
247 four methods of analysis on body and ear score, where each method considered different aspects of
248 the study design: (i) ignored the study design; (ii) considered correlation in the repeated
249 measurements; (iii) considered random effects from the hierarchical structure; (iv) considered the
250 correlation structure and the random effects. Table 1 provides a comparison of the different inferential
251 methods considered in this paper. Depending on the study design, it indicates which inferential
252 method would be appropriate for different types of data.

253 *(i) Ignoring study design (without accounting for repeated measures or hierarchical structure)*

254 To demonstrate the effects of ignoring the study design completely, i.e. not accounting for repeated
255 measures of individuals and random effects, we fitted a generalised linear model (GLM) to body and

256 ear score. Specifically a log linear model (LLM) was fitted to body score and a cumulative logistic
257 regression model (CLM) was fitted to ear score.

258 *(ii) Repeated measures (without accounting for hierarchical structure)*

259 As we assumed body score is continuous, we performed a multivariate analysis of covariance
260 (MANCOVA) with a Gaussian distribution. This methodology compares the means of all the different
261 possible groups and determines whether a significant difference is present when accounting for a
262 possible time-dependent correlation between the assessments. We accounted for the replications
263 within this inferential analysis using an error structure for individuals within replications.
264 MANCOVA assumes that the assessments measured are taken at equally spaced points in time, and
265 the difference in time is the same for each individual. Only individuals with complete data are
266 included.

267 As ear score is an ordinal variable, we fitted a cumulative logistic regression model for repeated
268 measures. To account for repeated measurements of the ear score, the parameters were estimated via
269 generalized estimating equations (GEE), which allow for the presence of a possible time-dependent
270 correlation between ear score assessments made at different times. However, a covariate for the
271 replication was also included to account for the possible differences between replications.

272 *(iii) Hierarchical structure (without accounting for repeated measures)*

273 To remove the effect of the repeated measures we produced a summary variable for each pig. The
274 summary variable for body score was simply the mean of the log transformed body score across each
275 of the three repeated measures. The summary variable for ear score was slightly more complicated.
276 Often categorical variables are summarised by their median or modal value. However, as the median
277 and mode are not influenced by extreme values, it meant that severe injuries were missed. Therefore,
278 we summed the ear score for each replication, then combined some of the categories according to the
279 frequency and level of injury the category represented to bring the score in line with the original
280 scoring system. The new ear score categories were 0 = less than 2 occurrences of superficial lesions,

281 or 1 occurrence of a deep lesion; 1 = 1 occurrence of a deep lesion and 1 occurrence of a superficial
282 lesion or 3 occurrences of superficial lesion; 2 = more than 1 occurrence of a deep lesion.

283 To account for the random effects of pen within replication we fitted a mixed effects linear regression
284 model (LME) to the mean log body score

$$285 \quad y_{i,j} = \alpha + X_{i,j}\boldsymbol{\beta} + Z_{i,j}\boldsymbol{\delta}_i,$$

286 **Equation 1**

287 and a cumulative logistic mixed effects regression model (CLME) to the re-categorized sum of ear
288 score

$$289 \quad \text{logit}(\text{Pr}[Y_{i,j} < k]) = \alpha_k + X_{i,j}\boldsymbol{\beta} + Z_{i,j}\boldsymbol{\delta}_i,$$

290 **Equation 2**

291 where: $y_{i,j}$ is the mean log body score; $Y_{i,j}$ is the ear score category for $k=0,1,2$; α is the intercept
292 whereas α_k is the intercept for the k^{th} cumulative logit; $\boldsymbol{\beta}$ is a vector of fixed effects coefficient
293 estimates; $X_{i,j}$ are the fixed covariates design vector for the j^{th} pig, in the i^{th} replication $\boldsymbol{\delta}_i$ is a vector of
294 the random effects for replication i ; and $Z_{i,j}$ is a design vector of the random effects.

295 An important difference between the GLM and a mixed effects model comes from the estimation of
296 the variance. In a GLM only the variance of the individual pigs is required, whereas now an estimate
297 for the variance for the individual pigs and the replications is required.

298 *(iv) Hierarchical data with repeated measures*

299 To account for both the hierarchical design and repeated measurements within this study, we fitted the
300 log linear and cumulative logistic, mixed effects model as defined in eEquation 3Equation 4:

$$301 \quad \log(y_{i,j,t}) = \alpha + X_{i,j,t}\boldsymbol{\beta} + Z_{i,j,t}\boldsymbol{\delta}_i,$$

302

Equation 3

303

$$\text{logit}(\Pr[Y_{i,j,t} < k]) = \alpha_k + X_{i,j,t}\boldsymbol{\beta} + Z_{i,j,t}\boldsymbol{\delta}_i.$$

304

Equation 4

305 These are very similar to Equation 1Equation 2, and in fact, the mathematical representation only
306 requires the addition of a subscript t to denote the time element in the random effects model. See
307 Twisk (2012) for more details on this type of analysis.

308 Computationally, as we are treating body score as a continuous Gaussian distributed variable,
309 estimation of the coefficients and the variance for the replications and individuals in Equation 3 can
310 be accomplished via GEE. However, there is no software available currently which can produce a
311 mixed effects cumulative logistic regression model with repeated measures where the correlation
312 between each observation depends on the time difference between repeated measures.). We concluded
313 that as we only had three repeated observations, estimation of the random effects was more important
314 than using GEE to account for a time dependent correlation structure for ear score. However, a
315 random effect term for each pig was included instead, as it assumes the correlation between
316 observations is constant over time.

317 *Small Sample Sizes*

318 To investigate the effects of small sample sizes, a repeated measures model was fitted to the data of
319 each replication. This led to 12 statistical models, one for each replication, which each consisted of 72
320 pigs per model/replication (18 pigs assigned to 1 of 4 pens), with a maximum of three skin lesion
321 assessments each, giving a total of number of observations of 216 per model. Each GLM consisted of
322 the same covariates, which were equivalent to the covariates in the final hierarchical repeated
323 measures model.

324

325

326 3. RESULTS

327 For 862 individual pigs we had a measurement for at least one of the injury assessments. For body
328 score there were two pigs with missing data for the first observation, seven pigs with missing data for
329 the second observation and nine pigs with missing data for the third observation. For ear score there
330 were three pigs with missing data for the first observation, seven pigs with missing data for the second
331 observation and 10 pigs with missing data for the third observation.

332 3.1 Body Score

333 3.1.1 Exploratory Analysis

334 The plots of the kernel smooth estimators in figure 2 a) – e) depict a cubic relationship with time. The
335 kernel estimators of log body score are between 1 and 2 at the first examination (day 0), with a
336 decline in log body score by the second examination (days 8-17), but by the third examination (days
337 29-39) there is an increase. All covariate groups mirror this pattern.

338 However, the slopes for each replication varied, as shown in figure 2 a), thus implying a random slope
339 for replication over time was required. Figure 2 b) of the Gaussian kernel smooth estimators for each
340 pen was used to determine whether different housing features were worth investigating. It is clear that
341 pigs within pen 3 tended to have a higher body score than any of the other three pens, which all
342 appeared to be quite similar. There was a difference between the intercept and a slight difference
343 between the slopes for each pen.

344 The plots in figure 2 c) to e) further identify differences between the pens. Comparing the score of the
345 different environments in figure 2 c), the difference between the less and more enriched environments
346 is only evident after approximately 14 days. This implies an interaction between time and
347 environment. The plot in figure 2 d) shows that pigs in the pens to the front of the experimental room
348 had a consistently higher body score than pigs in the pens located at the back. We also observed that
349 pigs in pens on the right side of the room had a higher body score than those in pens on the left side of
350 the room, as shown in figure 2 e).

351 The plot in figure 2 f) is a scatter plot of body score by standardised relative weight. The blue line is
352 the kernel smooth estimator using a bandwidth of 0.75. Less than 3% of the standardised weight
353 values were either > 2 or < -2 , which meant there were insufficient values to produce a reliable
354 estimate of the relationship between body score and relative weight. However, the plot suggested that
355 for a relative weight between -2 and 2 , the relationship was linear and as weight increased so did log
356 body score.

357 *3.1.2 Inferential Analysis*

358 Table 2 contains all the summary statistics for the fixed effects (coefficient estimate, standard error,
359 Student's t-value and p-value) for the most appropriate model, (iv) LLME + GEE, and the p-values
360 for all fixed effects for the three comparison methods, (i) LLM, (ii) MANCOVA and (iii) LLME. If a
361 p-value was greater than 0.05 it was not included in the table. In all the statistical models the
362 enrichment level, location of the pen (left/right side, front/back of the experimental room) was
363 significantly associated with body score. Relative weight was a significant component in 3 out of the
364 4 statistical models.

365 The LLME + GEE model accounted for a random intercept and slopes over time for pens within
366 replications, and a Gaussian correlation structure between observations for each pig. There was a
367 significant cubic relationship with time, this can also be seen in figure 2 (a)-(e) of the kernel
368 estimators. The significant relative weight coefficient implied that a unit increase in relative weight
369 resulted in a 0.05 increase in log body score, which equates to a 5% increase in body score. On
370 average, pigs on the right side of the room had a 0.094 higher log body score, i.e. their body score was
371 9.9% higher than those on the left side of the room. Also pigs with more enrichment and those in pens
372 located at the front of the experimental room had higher log body scores by 0.124 (13.2% increase in
373 body score) and 0.09 (9.4% increase in body score), respectively.

374

375 3.1.3 *Small Sample Sizes*

376 Figure 3 a) is a box plot of the coefficient estimate when using GEE to analyse each replication; when
377 the random effect for replication was not included, with the fixed effect coefficient estimates under
378 LLME + GEE model (table 2) included as a red cross. The box plot for relative weight was the only
379 one where the whiskers of the plot did not include zero, implying this was the only covariate with a
380 significant association with log body score for all but one replicate. This suggested that the coefficient
381 estimate for relative weight should remain fairly consistent across replications. For pen location (left/
382 right, front/back of the experimental room), and more enriched pens, the coefficient estimates showed
383 greater variance.

384 The median coefficient estimates were: weight = 0.04; right side of experimental room = 0.1; location
385 to the front = 0.14; and more enriched environment = 0.11. Comparing these values with the
386 coefficients estimates of the LLME + GEE model in table 2 we see that these values are quite similar,
387 and encouraging as a form of sensitivity analysis. Within one replication, there are 216 observations.
388 If we were to perform a t-test on these 216 observations to detect the largest effect size of 0.14 in log
389 body score, assuming the standard deviation was 0.6 (estimated from the entire dataset), then we
390 would have $\approx 40\%$ power to detect this difference. This does not account for the repeated measures,
391 which would reduce the power further.

392 3.2 *Ear Score*

393 3.2.1 *Exploratory Analysis*

394 From figure 4 there is evidence of a cubic relationship between ear score and time when comparing
395 the proportion of pigs with an ear score of 0 with 1 and/or 2 (all plots on the left), where there is a
396 decrease, plateau, then further decrease. However, the plots comparing the proportions observed in 0
397 and/or 1 with 2 (plots on the right) appear to be exponentially decaying.

398 The plots in figure 4 show the proportional change in the pigs observed within each ear score group
399 with Gaussian kernel estimators to convey how the relationship between ear score changes over time

400 for different housing features. In figure 4 a) the variability in the shape of the relationship between ear
401 score and time for the different replications indicate a different slope for each replication over time is
402 required. However, in figure 4 b) the estimators for each pen have a similar shape, but different
403 intercepts. There are clear differences in figures 4 c) and d) between environment and location next to
404 the front or the back of the experimental room.

405 *3.2.2 Inferential Analysis*

406 Table 3 shows all the summary statistics for fixed effects (coefficient estimate, standard error,
407 Student's t-value and p-value) for the cumulative logistic mixed effects regression model with random
408 effect for pigs, (iv) CLME +1, and significant p-values for fixed effects from the three comparator
409 methods (i) CLM, (ii) GEE and (iii) CLME. Within each statistical model, ear score was shown to
410 have a significant association with the level of enrichment and the front/back pen location.

411 The CLME+1 model included random intercept and slope terms for pen within replication to account
412 for the differences between replications over time, and a random intercept for each pig to account for
413 the correlation between repeated measures. To discuss our findings, we use odds ratios (i.e.
414 exponential transformation of the coefficients), so we can quantify the percentage increase or decrease
415 in odds that will result in the increase or decrease in ear injury score. In the CLME +1 model, pigs in
416 more enriched pens had 40% lower odds (Confidence Interval, CI: 14%, 58%) of having a higher ear
417 score compared to pigs in less enriched pens. Similarly, pigs in a pen located at the front of the room
418 had 33% lower odds (CI: 5%, 53%) of having a higher ear score.

419 *3.2.3 Small sample sizes*

420 We fitted a CLME model to each replication with a random intercept for each individual. Figure 3 b)
421 contains the box plot of the coefficient estimates from the ordinal logistic regression of ear score for
422 each replication. The fixed effect coefficient estimates under CLME+1 (table 3) are included as a red
423 cross in figure 3 b). There was a wide range of values for the coefficients from each replication
424 (median coefficient estimate for more enriched environment = -0.55; front of experimental room = -

425 0.21). Comparing the coefficient estimates for CLME and CLME+1, there was little difference
426 between pen enrichment estimates (0.04), but a larger difference between pen location estimates
427 (0.19).

428 ***3.3 Inference method comparisons***

429 For both types of injury score, the key associations between the injury score and environmental
430 factors were statistically significant across all four statistical models. Although, the magnitude of the
431 relationship and the direction was not always the same between the most appropriate statistical model
432 from approach (iv), and the other three statistical models, using methods (i) to (iii). The model via
433 approach (iii) for both injury scores provided no insight into changes in injury over time, as this
434 information was removed when summarising the injury scores.

435 Table 2 details the level of association between body score and the environmental factors for each
436 inferential method. Approach (i), the LLM, did not account for the repeated measure correlation or
437 random effects, and there was an additional significant association between body score and tail injury.
438 Whereas for approach (ii), the MANCOVA, which only accounted for repeated measurements, there
439 was a significant association between body score and sex. Neither of these associations were evident
440 in the exploratory analysis or in the most appropriate approach (iv). However, the association between
441 body score and weight was not statistically significant in approach (iii), the LLME model, but the
442 evidence from exploratory analysis and most appropriate model indicated there was a relationship
443 between these two variables.

444 In table 3 the statistical models from methods (i), CLM, and (ii), GEE, did not account for the random
445 effects of pen within replication that led to high order degree polynomials with the day, 7 and 5
446 respectively. There was no evidence in the exploratory analysis or the final most appropriate model
447 (CLME + 1), that this type of association between ear score and time was valid.

448

449

450 4. DISCUSSION

451 Comparing models where each incorporated different aspects of the study design demonstrated how
452 important using the most appropriate inferential analysis is when producing valid results. By
453 appropriately accounting for all sources of variation within the multilevel structure of the data (i.e.
454 pens within replications) and considering the potential time-dependent correlation between
455 observations, we increased the likelihood of identifying the true associations between the covariates
456 and injury scores. We also found that there was a strong agreement between exploratory and
457 inferential analysis, and associations seemed to be plausible.

458 In the most appropriate model for the data (repeated measures, mixed model), the strong significant
459 association of ear and body injury score with the non-linear time component is suggestive of a
460 complex relationship between behaviour and time. This observation was only possible because of the
461 repeated observations within pigs, and further validated by the replications of the study. Although the
462 variation in the inter-assessment interval time increased the statistical difficulty of the analysis, it did
463 mean that there was more information available about changes in injury score over a wider range of
464 interval differences. Ear and body injury score were both associated with the enrichment level and
465 front location of pen within the experimental room, although the direction of this association changed
466 for both covariates between injury scores. More enriched pens (coef. est. = -0.51, $p = 0.005$) and pens
467 at the front of the experimental room (coef. est. = -0.4, $p = 0.02$) were both associated with a
468 reduction in ear score, whereas those in more enriched pens (coef. est. = 0.09, $p = 0.02$), and pens at
469 the front of the experimental room (coef. est. = 0.11, $p = 0.003$) had a higher body score. Body score
470 was also associated with weight and pen location on the right side of the experimental room, such that
471 as weight increased so did body score (coef. est. = 0.05, $p < 0.001$), and those pigs in pens on the right
472 side of the experimental room also had a higher body score (coef. est. = 0.08, $p = 0.03$).

473 In this study, we investigated the impact of fitting statistical models that account for none, some and
474 all of the known structural features of a multilevel dataset. We also analysed the effect of small
475 sample size upon the most appropriate model. Similar investigations comparing inferential analyses

476 have been conducted in human and non-human medical literature (Hu et al., 1998; Wang and
477 Goonewardene, 2004), though this is the first example to the authors' knowledge in animal welfare.

478 In using an analytical approach that did not match the study design (approach (i): CLM), variance
479 within the dataset that was associated with either the hierarchical structure or the correlational
480 structure between repeated observations was not accounted for. This approach (CLM) led to
481 predictions of a complicated relationship between ear injury score and time, with a 7-degree
482 polynomial predicted to describe the relationship. For body score, the CLM predicted a cubic (i.e. 3-
483 degree polynomial) relationship with time, just as was predicted by the most appropriate model
484 (CLME+1). The high degree polynomial relationships predicted here result from poor estimation of
485 variance, due to the models attempting to explain variation in the data using only the covariates,
486 without the underlying hierarchical structure accounted for.

487 Including the correlation of the repeated measurements for approach (ii) via MANCOVA for body
488 score and GEE for ear score did increase the p-values, but it did not account for the substantial
489 variation caused by the random effects. Hence, there was an additional relationship between body
490 score and sex, and the association between ear score and day was now a 5-degree polynomial. One
491 substantial drawback back with MANCOVA is the strict format required of the data, i.e. equally
492 spaced repeated measures with no missing values. Using GEE analysis is more flexible and the
493 observations do not necessarily have to be equally spaced. However as the correlation coefficients
494 between repeated measurements of ear score were all less than 0.3, and the differences between the
495 estimators for replications and pens from the plots in figure 3 a) and b) appeared quite high, this
496 suggested the random effects terms for replication and pen were more important than accounting for
497 the correlation structure between repeated measurements. By replicating the study, we were able to
498 gain insight into differences between pens, which we had not considered for inclusion in our
499 experimental design prior to conducting the study; in particular, this would have been beneficial for
500 the location of the pens within the experimental room. Although we accounted for differences in
501 noise level with left/right side counter-balancing of the treatments, and accounted for potential

502 differences between pens at the front (near the door) versus at the back of the room with front/back
503 counter-balancing of treatments, we did not rotate the pens, which would have allowed us to account
504 for the additional locational differences detected in the data. Although we were unable to fully explain
505 the reason for differences between pen locations within the experimental room, we were able to
506 identify that pen location was a source of variation and we could therefore statistically remove any
507 undue influence this was having on other covariates within the model. Differences observed between
508 replications could be related to weather conditions, handlers and many other features not measured as
509 part of this study. Despite being unable to quantify all variation between replications, we believe that
510 replication on other farm sites would help to build up a more general picture across contexts.

511 Summary measures of both body and ear score were used in approach (iii), which resulted in lost
512 information about the nature of the relationships of body and ear score across time. Using this
513 approach, we were unable to identify a significant association between body score and weight via the
514 LLME model, but we detected a significant relationship between ear score and weight using the
515 CLME, as compared to the final appropriate model.

516 In the final approach (iv) for body score and ear score, there was evidence of a cubic relationship with
517 time for both injury scores. However, the direction of the coefficient estimates for day, day² and day³
518 differed between body and ear injury scores. For body scores, the coefficients for time were positive
519 for day and day² and negative for day³, whereas for ear score they were negative for day and day³ and
520 positive for day². This result implies that the underlying behaviour indicated by proxy from these
521 injury scores changed over time. For example, the initial decline in scores could be associated with
522 pigs becoming acquainted with one another as a hierarchy within a pen was established within the
523 first week (Barnett et al., 1994; Arey, 1999).

524 In both the final ear score and body score statistical models there was a significant association with
525 pen location (front/back of the room) and enrichment level (see section 3.2.2). Pigs in pens located at
526 the front of the room had lower odds of having a higher ear score (table 3), but higher odds of a higher
527 body score (table 2). Pigs in more enriched pens had lower ear scores (as described in section 3.2.2,

528 table 3). This result supports previous findings that aggressive events are reduced in larger pen sizes
529 (Fraser et al., 1991; Turner et al., 2000). Whereas the LME + GEE model for body score implies that
530 more enriched pens resulted in higher body injury scores.

531 Finding clear differences in the predictors for ear and body scores lends support to the hypothesis that
532 they have different underlying causes. Injuries to the ear are mainly received during aggressive
533 interactions (McGlone, 1985). Injuries to the body on the other hand, whilst accrued through
534 aggression, can also be the result of increased activity and play (Munsterhjelm et al. 2009; Camerlink
535 et al., 2013). Unfortunately, as tails were docked at birth we were not able to use tail injury as another
536 comparator, although research suggests that the majority of tail injuries reflect exploratory motivation
537 rather than aggression (Taylor et al., 2010). Applying a similar study to undocked pigs may provide
538 further detailed insight into aggression and the underlying motivating behaviours that lead to injuries.

539 Statistical techniques used to determine the validity in medical screening tests, such as a receiver
540 operator curve (ROC) analysis (Fawcett, 2006) or Bland-Altman test (Bland & Altman, 1986), may be
541 used to compare indicators of aggression to determine if they are a measure of the same quantity.

542 Whilst the final model selected is appropriate for the experimental design, it is not perfect. There are
543 currently no developed statistical methods available to analyse categorical outcome variables with a
544 time dependent correlation structure between repeated measures within a hierarchical model (such as
545 the random effects of replications within pens described within section 2.1). As such, we could not
546 account for both the correlational structure and hierarchy of the study design within current statistical
547 methodology. One possible solution could be to develop a statistical model with a probit link rather
548 than a logit link, as the probit link is associated with the Gaussian distribution, and it may be easier to
549 define a time dependent correlation structure with this compared to the logit link. However, the
550 interpretation of the probit link can be difficult as there are no direct interpretations of the coefficients,
551 instead it is necessary to refer to the marginal effects of the regressors (see Liao (1994) for more
552 details), and the estimation of the coefficients would be computationally intensive.

553 Differences between the results of the four inferential methods highlight the importance of initial
554 exploratory analysis in determining whether resulting significant associations are realistic, particularly
555 as all four methods used are technically appropriate, albeit with varying degrees of fit to the
556 experimental design. Strong evidence of a relationship in the exploratory analysis should translate to a
557 significant association observed within the inferential analysis. Although measures were taken into
558 account for layout of the experimental room, it was not possible to completely account for the extent
559 of this effect, and it was through exploratory analysis that we were provided with greater insight into
560 the magnitude and nature of the effect.

561 By analysing each replication separately, we were able to demonstrate how sample size affects the
562 final coefficient estimates. The decrease in data resulted in insufficient power to detect significant
563 associations, although the calculated medians of almost all the replications' coefficient estimates were
564 consistent with our full final models. The results clearly demonstrate that analysis of small sample
565 sizes may lead investigators to believe there was no association between the indicators for aggression
566 and covariates, whereas it could be the study is under-powered to detect the effect size (i.e. the
567 conclusion would be a type 2 error). As a simple demonstration, we performed a power calculation to
568 detect a mean difference in body score of 0.18 and standard deviation of 0.6, based on summary
569 statistics of enrichment level in the fifth week. The power calculation found that to detect such a
570 difference with 80% power at the 5% level of significance, a sample size of 176 pigs (total 352)
571 assigned to each enrichment level was required.

572 This study demonstrates through examples, how the type of indicator measured, the sample size and
573 choice of statistical analysis can affect model outputs and conclusions drawn. We also highlight the
574 importance of using an appropriate indicator to reflect the behaviour under investigation. The correct
575 inferential analysis is important for meaningful results, which are not only plausible, but also
576 supported by the exploratory analysis. To ensure the quality of animal science reports it is vital that a
577 study consists of an appropriate sample size, with statistical analysis appropriate for the study design.
578 These findings provide further support for the ARRIVE guidelines, but we feel that additional steps

579 may improve the quality of research by ensuring studies are designed based upon the inferential
580 analysis best equipped to answer the research question. It may be valuable to consider following
581 similar procedures as in medical trials with the formulation of a protocol and detailed documentation
582 of any unexpected and additionally planned deviations, which may subsequently affect the inferential
583 analysis. This way, while best laid plans may still go awry in practice, there will be a clear plan to
584 ensure that robust and appropriate analysis of the data can still be conducted.

585 **ACKNOWLEDGEMENTS**

586 The authors would like to thank AFBI for use of experimental room and care of the animals.

587 **ETHICS STATEMENT**

588 All procedures described were approved by the University of Lincoln's Ethics Committee on
589 8/9/2015, code COSREC62. This research was conducted at the Agri-Food and Biosciences Institute,
590 Northern Ireland and conformed to the Association for the Study of Animal Behaviour's guidelines on
591 the use of animals in research: <http://asab.nottingham.ac.uk/ethics/guidelines.php>.

592 **FUNDING STATEMENT**

593 This work was supported by the BBSRC (grants BB/K002554/1 and BB/K002554/2). MF was
594 supported by a Department for Employment and Learning Northern Ireland studentship and Queen's
595 University Belfast.

596 **REFERENCES**

- 597 Andersen, I.L., Andenæs, H., Bøe, K.E., Jensen, P., Bakken, M., 2000. The effects of weight
598 asymmetry and resource distribution on aggression in groups of unacquainted pigs. *Appl. Anim.*
599 *Behav. Sci.* 68, 107-120.
- 600 Arey, D.S., 1999. Time course for the formation and disruption of social organisation in group-housed
601 sows. *Appl. Anim. Behav. Sci.* 62, 199-207.

- 602 Arey, D.S., Franklin, M.F., 1995. Effects of straw and unfamiliarity on fighting between newly mixed
603 growing pigs. *Appl. Anim. Behav. Sci.* 45, 23-30.
- 604 Barnett, J.L., Cronin, G.M., McCallum, T.H., Newman, E.A., 1994. Effects of food and time of day
605 on aggression when grouping unfamiliar adult pigs. *Appl. Anim. Behav. Sci.* 39, 339-347.
- 606 Bates, D., Maechler, M., Bolker, B., Walker, S., 2015. lme4: Linear mixed-effects models using Eigen
607 and S4, <http://CRAN.R-project.org/package=lme4>.
- 608 Beattie, V.E., Walker, N., Sneddon, I.A., 1996. An investigation of the effect of environmental
609 enrichment and space allowance on the behaviour and production of growing pigs. *Appl. Anim.
610 Behav. Sci.* 48, 151-158.
- 611 Biau, D.J., Kernéis, S., Porcher, R., 2008. Statistics in Brief: The Importance of Sample Size in the
612 Planning and Interpretation of Medical Research. *Clinical Orthopaedics and Related Research* 466,
613 2282-2288.
- 614 Bland, J. M. & Altman, D. G. (1986). Statistical methods for assessing agreement between two
615 methods of clinical measurement. *Lancet* i, 301-310.
- 616 Borenstein, M., Hedges, L. V., Higgins, J. P. T., & Rothstein, H. R. (2009). Introduction to meta-
617 analysis. Chichester: Wiley.
- 618 Box, G.E., Draper, N.R., 1987. Empirical model-building and response surfaces. Wiley New York.
- 619 Brown, A.W., Mehta, T.S., Allison, D.B., 2017. Publication bias in science: What is it, why is it
620 problematic, and how can it be addressed? In: *The Oxford Handbook of the Science of Science
621 Communication*. Ed: K. Hall Jamieson, D. Kahan, D.A. Scheufele. Oxford University Press,
622 New York, USA.
- 623 Brunberg, E., Wallenbeck, A., Keeling, L.J., 2011. Tail biting in fattening pigs: Associations between
624 frequency of tail biting and other abnormal behaviours. *Appl. Anim. Behav. Sci.* 133, 18-25.
- 625 Camerlink, I., Turner, S. P., Bijma, P., Bolhuis, J. E. (2013). Indirect genetic effects and housing
626 conditions in relation to aggressive behaviour in pigs. *PloS one.*;8:e65136.
- 627 Carr, J., 1998. Garth Pig Stockmanship Standards. 5m Publishing.

- 628 Christensen, R.H.B., 2015. ordinal: Regression Models for Ordinal Data, [http://www.cran.r-](http://www.cran.r-project.org/package=ordinal)
629 [project.org/package=ordinal](http://www.cran.r-project.org/package=ordinal).
- 630 Conte, S., P. G. Lawlor, N. O'Connell, and L. A. Boyle. 2012. Effect of split marketing on the
631 welfare, performance, and carcass traits of finishing pigs¹. *J. Anim. Sci.* 90:373-380.
- 632 Cronbach, L. J. & Meehl, P. E. (1955). Construct validity in psychological tests. *Psychological*
633 *Buletin.* 52: 281-302.
- 634 Desire, S., Turner, S.P., D'Eath, R.B., Doeschl-Wilson, A.B., Lewis, C.R.G., Roehe, R., 2016.
635 Prediction of reduction in aggressive behaviour of growing pigs using skin lesion traits as selection
636 criteria. *Animal* 10, 1243-1253.
- 637 Drickamer, L.C., Arthur, R.D., Rosenthal, T.L., 1999. Predictors of social dominance and aggression
638 in gilts. *Appl. Anim. Behav. Sci.* 63, 121-129.
- 639 Fawcett, T. (2006). An introduction to ROC analysis. *Pattern Recognition Letters* 27: 861-874.
- 640 Fraser, D., Phillips, P.A., Thompson, B.K., Tennessen, T., 1991. Effect of straw on the behaviour of
641 growing pigs. *Appl Anim Behav Sci* 30, 307-318.
- 642 Freiman , J.A., Chalmers , T.C., Smith , H.J., Kuebler , R.R., 1978. The Importance of Beta, the Type
643 II Error and Sample Size in the Design and Interpretation of the Randomized Control Trial. *New*
644 *England Journal of Medicine* 299, 690-694.
- 645 Gulin, J., Rocco, D., Garca-Bournissen, F., 2015. Quality of Reporting and Adherence to ARRIVE
646 Guidelines in Animal Studies for Chagas Disease Preclinical Drug Research: A Systematic
647 Review. *PLoS Negl. Trop. Dis.* 9, 1-17.
- 648 Hopewell, S., Loudon, K., Clarke, M. J., Oxman, A. D., Dickersin, K. (2009). Publication bias in
649 clinical trials due to statistical significance or direction of trial results. *Cochrane Database of*
650 *Systematic Reviews*, Issue 1. Art. No.: MR000006.
- 651 Hu, F.B., Goldberg, J., Hedeker, D., Flay, B.R., Pentz, M.A., 1998. Comparison of population-
652 averaged and subject-specific approaches for analyzing repeated binary outcomes. *American*
653 *Journal of Epidemiology* 147, 694-703.

- 654 Ioannidis, J. P., Allison, D. B., Ball, C. A., Coulibaly I., Cui X., Culhane, A. C., Falchi, M.,
655 Furlanello, C., Game, L., Jurman, G., Mangion, J., Mehta, T., Nitzberg, M., Page, G. P., Petretto,
656 E., van Noort, V. (2009). Repeatability of published microarray gene expression analyses. *Nature*
657 *Genetics*. 41:149–155
- 658 Kilkenny, C., Browne, W.J., Cuthill, I.C., Emerson, M., Altman, D.G., 2010. Improving Bioscience
659 Research Reporting: The ARRIVE Guidelines for Reporting Animal Research. *PLoS Biol* 8, 1-5.
- 660 Kilkenny, C., Parsons, N., Kadyszewski, E., Festing, M.F.W., Cuthill, I.C., Fry, D., Hutton, J.,
661 Altman, D.G., 2009. Survey of the Quality of Experimental Design, Statistical Analysis and
662 Reporting of Research Using Animals. *PLoS ONE* 4, 1-11.
- 663 Lahrmann, H.P., Oxholm, L.C., Steinmetz, H., Nielsen, M.B.F., D'Eath, R.B., 2015. The effect of
664 long or chopped straw on pig behaviour. *Animal* 9, 862-870.
- 665 Liao, T.F., 1994. Interpreting probability models: Logit, probit, and other generalized linear models.
666 Sage.
- 667 McCance, I., 1995. Assessment of statistical procedures used in papers in the Australian Veterinary
668 Journal. *Aus. Vet. J.* 72, 322-329.
- 669 McGlone, J.J., 1985. A Quantitative Ethogram of Aggressive and Submissive Behaviors in Recently
670 Regrouped Pigs. *J. Anim. Sci.* 61, 556-566.
- 671 Morgan, C.A., Deans, L.A., Lawrence, A.B., Nielsen, B.L., 1998. The effects of straw bedding on the
672 feeding and social behaviour of growing pigs fed by means of single-space feeders. *Appl. Anim.*
673 *Behav. Sci.* 58, 23-33.
- 674 Munsterhjelm, C., Peltoniemi, O. A., Heinonen, M., Halli, O., Karhapaa, M., et al. (2009). Experience
675 of moderate straw bedding affects behaviour of growing pigs. *Appl Anim Behav Sci.* 118:42–53.
- 676 Nettle, D., Monaghan, P., Boner, W., Gillespie, R., Bateson, M., 2013. Bottom of the Heap: Having
677 Heavier Competitors Accelerates Early-Life Telomere Loss in the European Starling, *Sturnus*
678 *vulgaris*. *PLoS ONE* 8, e83617.

- 679 Sargeant, J.M., Thompson, A., Valcour, J., Elgie, R., Saint-Onge, J., Marcynuk, P., Snedeker, K.,
680 2010. Quality of Reporting of Clinical Trials of Dogs and Cats and Associations with Treatment
681 Effects. *J. Vet. Intern. Med.* 24, 44-50.
- 682 Statham, P., Green, L., Mendl, M., 2011. A longitudinal study of the effects of providing straw at
683 different stages of life on tail-biting and other behaviour in commercially housed pigs. *Appl.*
684 *Anim. Behav. Sci.* 134, 100-108.
- 685 Taylor, N.R., Main, D.C.J., Mendl, M., Edwards, S.A., 2010. Tail-biting A new perspective.
686 *Veterinary Journal* 186, 137-147.
- 687 Team, R.C., 2015. R: A Language and Environment for Statistical Computing, R Foundation for
688 Statistical Computing, Vienna, Austria.
- 689 Touloumis, A., 2016. multgee: GEE Solver for Correlated Nominal or Ordinal Multinomial
690 Responses, <http://www.cran.r-project.org/package=multgee>.
- 691 Turner, S.P., Ewen, M., Rooke, J.A., Edwards, S.A., 2000. The effect of space allowance on
692 performance, aggression and immune competence of growing pigs housed on straw deep-litter at
693 different group sizes. *Livest. Prod. Sci* 66, 47-55.
- 694 Twisk, J.W.R., 2012. *Applied Longitudinal Data Analysis for Epidemiology: A Practical Guide*. 2 ed.
695 Cambridge University Press, Cambridge.
- 696 Vogt L, Reichlin TS, Nathues C, Würbel H (2016) Authorization of Animal Experiments Is Based on
697 Confidence Rather than Evidence of Scientific Rigor. *PLoS Biol* 14(12): e2000598.
- 698 Wand, M.P., Jones, M.C., 1994. *Kernel Smoothing*. Chapman & Hall/CRC.
- 699 Wang, Z., Goonewardene, L.A., 2004. The use of MIXED models in the analysis of animal
700 experiments with repeated measures data. *Canadian Journal of Animal Science* 84, 1-11.
- 701
- 702
- 703

704 **FIGURE LEGENDS**

705 **Figure 1:** The six-point scaling system used to assess injuries to pig's body areas and outline of body
706 areas for injury scoring; Ears, Snout, Shoulders, Legs, Back, Flanks, Hind quarters and Tail.

707 **Figure 2:** Plots of the log transformed body score by day with a Gaussian kernel smooth estimator
708 with a bandwidth of 15 for a) replication; b) pen; c) enrichment; d) location to the front or back of the
709 experimental room; e) location on either side of the experimental room. The light grey area depicts the
710 time period the second injury assessments were gathered, all points gathered after this period are the
711 third injury assessments and all points before are the first; f) Plot of the pig's relative weight for each
712 pen within replication by log body score with a Gaussian kernel smooth estimator with bandwidth of
713 4. The grey area of the plot indicates the region where 95% of the data is located, and where the
714 kernel estimator will be most reliable.

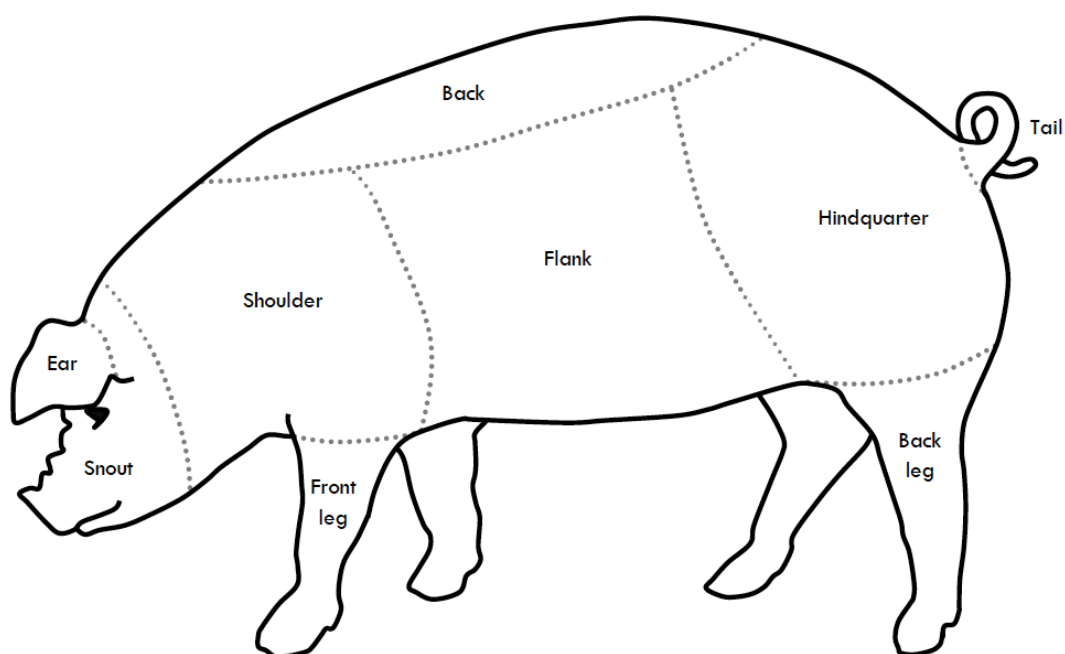
715 **Figure 3:** a) Box plot of the fixed effect coefficient estimates for the log linear regression model for
716 body score for each replication. The red crosses represent the fixed effect coefficient estimates for the
717 LLME + GEE from table 2. b) Box plot of the fixed coefficient estimates from the ordinal logistic
718 regression of ear score for each replication. The red crosses represent the fixed effect coefficient
719 estimates for the CLME +1 in table 3.cross.

720 **Figure 4:** Left plots: observed proportion with an ear score of 0 and 1/2. Right plots: observed
721 proportion with an ear score of 0/1 and 2, with Gaussian kernel estimators with a bandwidth of 15 for
722 a) replications; b) pens; c) enrichment; or d) location to the front or the back of the experimental
723 room. The light grey area depicts the time period the second injury assessments were gathered, all
724 injury assessments gathered after this period are the third injury assessments and all injury
725 assessments before are the first.

726

727

728 **Figure 1**



729

| Score | Scaling System |
|-------|--|
| 0 | No injuries. |
| 1 | One small superficial lesion. |
| 2 | More than one small, superficial lesion; or just one red (deeper than score 1) but still superficial lesion. |
| 3 | One or several big and deep lesions. If deep, only one single lesion. If not so deep, several red lesions. |
| 4 | One very big, deep and red lesion. Or many deep, red lesions. |
| 5 | Many, very big, deep and red lesions covering the skin area. |

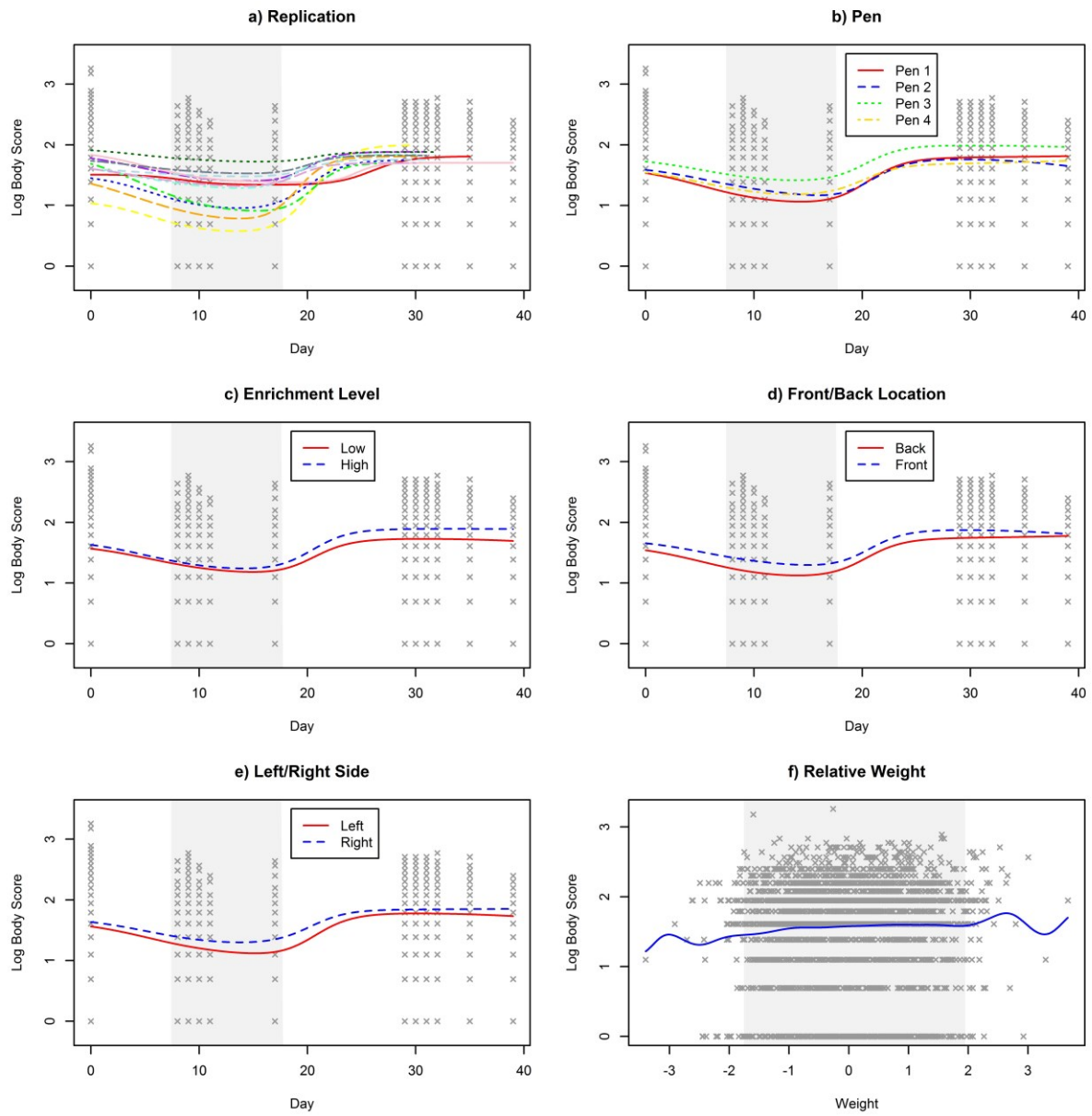
730

731

732

733 **Figure 2**

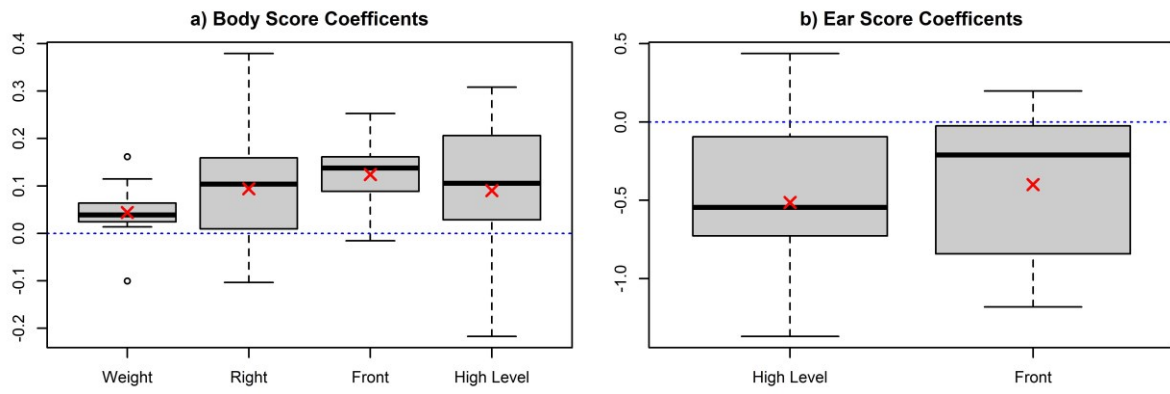
734



735

736

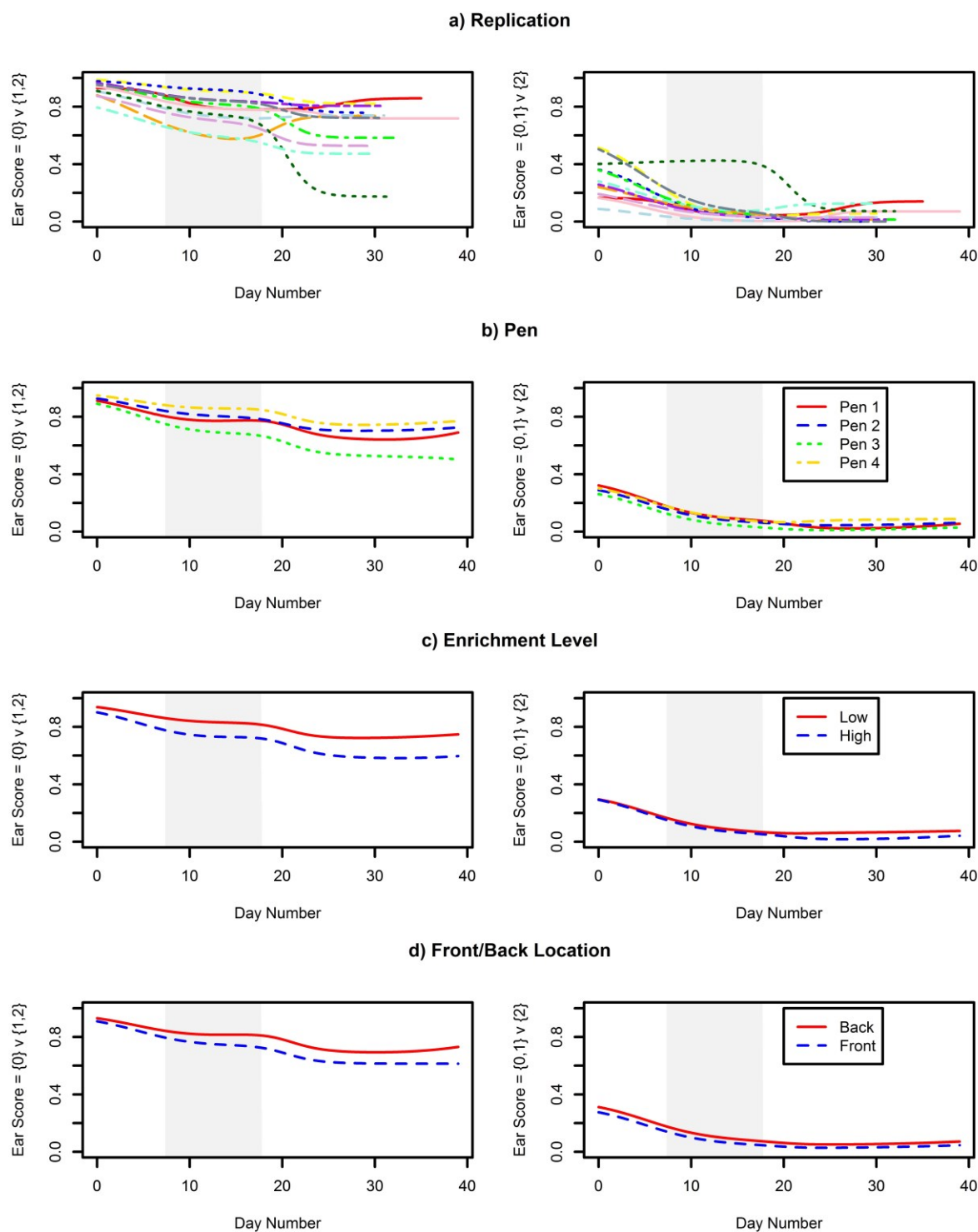
737 **Figure 3**



738

739

740 **Figure 4**



741

742

743 **Table 1**

| Data | Inferential Method | | | | |
|--------------------------------|---------------------------|------------|------------|------------|------------------|
| | <i>MANCOVA</i> | <i>GLM</i> | <i>LME</i> | <i>GEE</i> | <i>LME + GEE</i> |
| Univariate | | C O | | | |
| Multivariate | C | | | | |
| Repeated | | | | C O | |
| Hierarchical | | | C O | | |
| Repeated + Hierarchical | | | | | C |

744 **Table 1:** Types of data that can be analysed using different inference methods, where C represents
 745 continuous data and O represents ordinal data. MANCOVA=Multivariate Analysis of Covariance;
 746 GLM=Generalised linear model; LME=Linear mixed effects model; GEE=General Estimating
 747 Equation model.

748

749 **Table 2**

| | LLME + GEE | | | | LLME | MANCOVA | LLM |
|------------------------|------------|-----------|----------|----------|----------|---------|----------|
| | <i>n</i> | | | | <i>n</i> | | |
| Pigs | 862 | | | | 862 | 855 | 862 |
| Body Score | 2565 | | | | 862 | 2550 | 2556 |
| | β | <i>SE</i> | <i>t</i> | <i>p</i> | <i>p</i> | | |
| Day | 5.87 | 2.47 | 2.38 | 0.0173 | | | < 0.0001 |
| Day² | 11.45 | 2.35 | 4.87 | < 0.0001 | | | < 0.0001 |
| Day³ | -6.39 | 1.30 | -4.93 | < 0.0001 | | | < 0.0001 |
| More Enriched | 0.09 | 0.04 | 2.40 | 0.0224 | 0.0151 | 0.0003 | 0.0003 |
| Location: Right | 0.08 | 0.04 | 2.26 | 0.0307 | 0.0109 | 0.0018 | < 0.0001 |
| Sex | | | | | | 0.0041 | |
| Weight | 0.05 | 0.01 | 3.41 | 0.0007 | | 0.0278 | 0.0013 |
| Location: Front | 0.11 | 0.04 | 3.16 | 0.0034 | 0.0011 | 0.0003 | < 0.0001 |

750 **Table 2:** Summary statistics for inferential analysis of Body Score via the: log linear mixed effects
 751 model for repeated measures (LLME + GEE); linear mixed effects model of pig's mean log body
 752 score (LME); multivariate analysis of covariance (MANCOVA) of log body score, and a log linear
 753 regression model (LLM). Where: *n* is the number of pigs/body score assessment; β is the parameter
 754 estimate; *SE* is the standard error; *t* is the Student's t test statistic and *p* is the probability value
 755 associated with each covariate. **Day** is the day within the trial that observations were recorded; **More**
 756 **Enriched** refers to pens that had more enrichment (compared with Less Enriched); **Location: Right**
 757 refers to pens on the right side of the room (compared to pens on the left side of the room); **Location:**
 758 **Front** refers to pens at the front of the room (compared to pens at the back of the room).

759

760 **Table 3**

| | CLME + 1 | | | | CLME | GEE | CLM |
|------------------------|----------|-----------|----------|----------|----------|----------|----------|
| | <i>n</i> | | | | <i>n</i> | | |
| Pigs | 862 | | | | 862 | 862 | 862 |
| Ear Score | 2572 | | | | 862 | 2572 | 2572 |
| | β | <i>SE</i> | <i>t</i> | <i>p</i> | <i>p</i> | | |
| Day | -51.68 | 5.75 | -8.99 | < 0.0001 | | < 0.0001 | < 0.0001 |
| Day² | 31.30 | 5.74 | 5.45 | < 0.0001 | | < 0.0001 | < 0.0001 |
| Day³ | -13.56 | 6.51 | -2.08 | < 0.0369 | | 0.0453 | 0.0003 |
| Day⁴ | | | | | | < 0.0001 | < 0.0001 |
| Day⁵ | | | | | | 0.0194 | < 0.0001 |
| Day⁶ | | | | | | | 0.0255 |
| Day⁷ | | | | | | | < 0.0001 |
| More Enriched | -0.51 | 0.18 | -2.79 | 0.0053 | 0.0131 | < 0.0001 | < 0.0001 |
| Weight | | | | | 0.0302 | | |
| Location: Front | -0.40 | 0.18 | -2.25 | 0.0247 | 0.0328 | < 0.0001 | < 0.0001 |

761

762

763 **Table 3:** Summary statistics for inferential analysis of Ear Score via the: cumulative logistic mixed
764 effects model with rep, pen and pig random effects (CLME + 1); cumulative logistic mixed effects
765 model with rep and pen random effects for summary ear score (CLME); cumulative logistic
766 regression model for repeated measures (GEE); the cumulative logistic regression model (CLM).
767 Where: n is the number of pigs/ear score assessment; β is the parameter estimate; SE is the standard
768 error; t is the Student's t test statistic and p is the probability value associated with each covariate.
769 **Day** is the day within the trial that observations were recorded; **More Enriched** refers to pens that
770 had more enrichment (compared with Less Enriched); **Location: Front** refers to pens at the front of
771 the room (compared to pens at the back of the room).

772