# Ordinal or visual analogue scales for assessing aspects of broiler chicken welfare?

SCHOLARONE™
Manuscripts

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Contact dermatitis on the breast and abdominal areas (CD)

VAS (mm)

| 0 | 16 | 40 | 68 | 100 |

| 26.6% | 49.7% | 17.5% | 6.1% |

ORS

| 0 | 1 | 2 | 3 |

Footpad dermatitis (FP)

VAS (mm)

| 0 | 20 | 30 | 58 | 100 |

| 36.6% | 9.5% | 23.4% | 26.2% + 4.2% |

ORS

| 0 | 1 | 2 | 3+4 |

Hock burn (HB)

VAS (mm)

| 0 | 20 | 34 | 56 | 100 |

| 60.1% | 29.9% | 6.9% | 2.1% + 1.0% |

ORS

| 0 | 1 | 2 | 3+4 |

Bird soiling (BS)

VAS (mm)

| 0 | 40 | 70 | 100 |

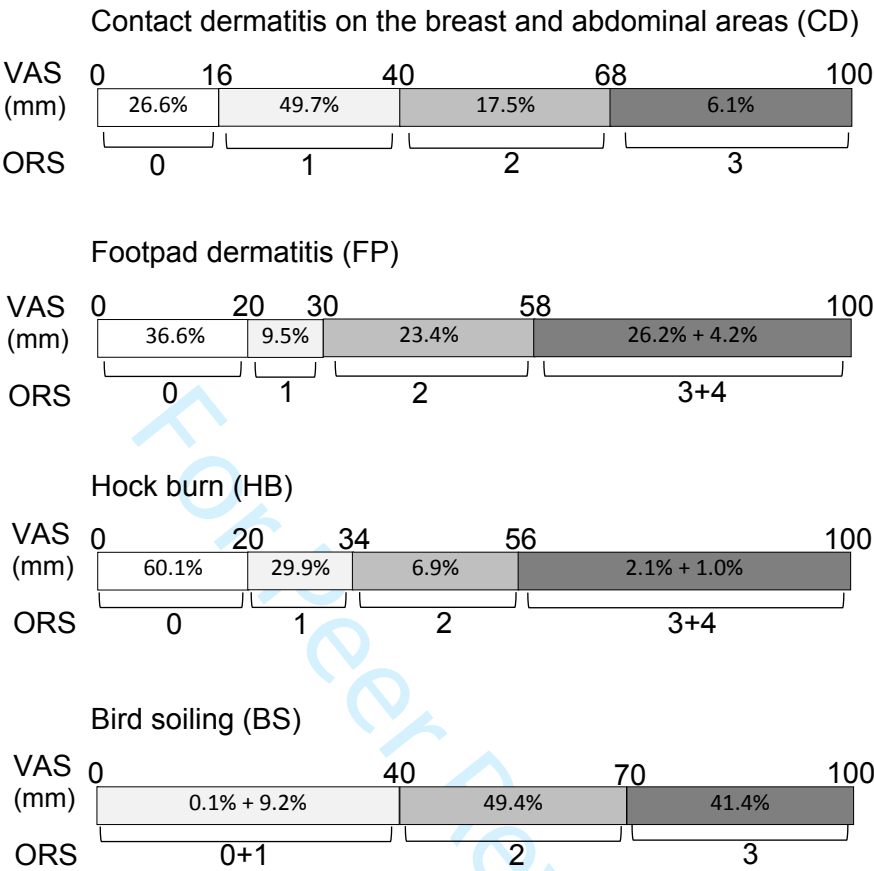| 0.1% + 9.2% | 49.4% | 41.4% |

ORS

| 0+1 | 2 | 3 |

Table 1. Estimates of inter-rater reliability and confidence interval, 5,000 bootstrap samples, for animal welfare indicators from 1,303 broiler chickens assessed on farm by three raters using both ordinal scale (ORS) and visual analogue scale (VAS).

| Welfare indicator | Scale | Intraclass correlation | Confidence interval (95%) | P-value* |
|---|---|---|---|---|
| Contact dermatitis on the breast and abdominal areas | ORS | 0.68 | (0.58 - 0.77) | <0.001 |
| | VAS | 0.77 | (0.67 - 0.85) | |
| Footpad dermatitis | ORS | 0.91 | (0.87 - 0.93) | <0.001 |
| | VAS | 0.88 | (0.83 - 0.92) | |
| Hock burns | ORS | 0.67 | (0.55 - 0.76) | <0.001 |
| | VAS | 0.72 | (0.60 - 0.80) | |
| Bird soiling | ORS | 0.61 | (0.46 - 0.73) | 0.447 |
| | VAS | 0.54 | (0.36 - 0.69) | |

Table 2a. Correlation of ordinal scale (ORS) and visual analogue scale (VAS) for the mean of values given by the three raters and for the individual values of each rater.

Table 2b. Correlation of broiler chicken welfare indicators measured on farm using ORS and VAS, 1,303 birds.

Table 2a

| Indicator | Spearman rank correlation between ORS and VAS* | |
|---|---|---|
| | Mean | Individual |
| Contact dermatitis on the breast and abdominal areas (CD) | 0.96 | 0.89 |
| Footpad dermatitis (FP) | 0.97 | 0.95 |
| Hock burn (HB) | 0.90 | 0.77 |
| Bird soiling (BS) | 0.94 | 0.81 |

*P < 0.0001

Table 2b

Correlation between indicators* (ORS, Spearman correlation; VAS, Pearson correlation)

| Indicator | Scale | FP | HB | BS |
|---|---|---|---|---|
| CD | ORS | 0.06 | 0.24 | 0.34 |
| | VAS | 0.09 | 0.35 | 0.34 |
| FP | ORS | | 0.17 | 0.08 |
| | VAS | | 0.26 | 0.12 |
| HB | ORS | | | 0.25 |
| | VAS | | | 0.24 |

1 **Ordinal or visual analogue scales for assessing aspects of broiler chicken welfare?**

2

3 **Abstract**

4 Information may be lost when the gradation of animal welfare is scored through ordinal

5 scales. Therefore, some advocate the use of continuous scales, which may be tagged with

6 internal anchors. Equidistant tags are used; however, studies have demonstrated that

7 empirical data for the space between tags tend to be non-equidistant. Ordinal rate scales

8 (ORS) and visual analogue scales (VAS) were tested for the assessment of contact

9 dermatitis on the breast and abdominal areas (CD), footpad dermatitis (FP), hock burns (HB)

10 and bird soiling (BS) in broiler chickens. Calculations regarding the inter-rater reliability, the

11 correlation between VAS and ORS and amongst the welfare indicators measured with both

12 scales, as well as the equidistance of ORS categories in relation to values measured using

13 VAS, were made. A total of 1,303 broiler chickens from 10 flocks was assessed on-farm by

14 three trained raters using an ORS and a VAS anchored only with the minimum and the

15 maximum scores at each end. Inter-rater reliabilities of CD (0.68 vs 0.77, P<0.001) and HB

16 (0.67 vs 0.72, P<0.001) were higher when using VAS compared with ORS, but that of FP

17 (0.91 vs 0.88, P<0.001) was lower. Correlations between ORS and VAS varied between

18 0.90-0.97 and 0.77-0.95 (P<0.001) respectively, considering mean and individual values of

19 the three raters. Low to moderate correlations were observed between the four indicators

20 using ORS and VAS. Tags on VAS that best represented ORS were non-equidistant.

21 Results suggest both scales were reliable to assess the selected broiler chicken welfare

22 indicators.

23

24 **Keywords**: animal welfare, animal-based measures, categorical scale, continuous scale,

25 poultry

26

1. **Introduction**

The development and application of protocols to assess animal welfare (AW) has increased worldwide. In addition to registering absence and presence of AW issues, it is often useful and informative to score gradations of these issues. Assuming equal reliability, the more refined these gradations are scored, the more sensitive becomes the detection of relevant AW aspects, such as AW progress over time, differences between the welfare of groups of animals or effects interventions have on the lives of animals. Scientific research has encouraged the development of new techniques to assess AW in field conditions. Reliability between raters is an important criterion in the selection of AW indicators, since there is high probability of single person assessments due to manpower costs of animal-based monitoring schemes (Tuyttens et al., 2014). There are some initiatives for assessing the welfare of broiler chickens, like the Welfare Quality® (2009), the AssureWel (2014) and the Global Animal Partnership® (2018). These protocols include measures, predominantly presented as ordinal rating scales (ORS) ranging from 2- to 6-point scales. Raters can be trained to score reliably using ORS, and much of advances in knowledge of broiler chicken welfare are due to the application of ORS in the assessment of welfare in experimental and commercial flocks.

Descriptors, photos and videos may be used for illustrating, and practicing the recognition of stepwise increases in severity, thereby increasing consistency within and between observers. This also implies that data from different studies can be compared if the same ORS are used. However, assessing continuous welfare traits by using discontinuous scales may be disadvantageous (Tuyttens et al., 2009). The use of ORS may result in reduced sensitivity when raters are able to discriminate more levels of the assessed indicator than the number of categories allow for and are forced to group gradations they perceive as different into the same category.

2

53    A different type of scale, the visual analogue scale (VAS), is largely used to assess

54    pain in humans and non-human animals (de Grauw and van Loon, 2016; Hjermstad et al.,

55    2011). In AW assessment, VAS has also been applied to assess qualitative behavior defined

56    by Wemelsfelder et al. (2001) as one of the "whole animal" measures that aim to assess the

57    overall subjective experience or mood of an animal (Fleming et al., 2016; Grosso et al.,

58    2016; Minero et al., 2016), and lameness (Flower and Weary, 2006; Nalon et al., 2014;

59    Tuyttens et al., 2009; Vieira et al., 2015) in different species. VAS is a continuous scoring

60    system that consists of a line, which varies usually from 100 to 125 mm in length, anchored

61    by the minimum and the maximum score at each end. Thus, VAS removes the constraint of

62    grouping information into discrete units and enables raters to achieve greater sensitivity in

63    their scoring for aspects that vary along a continuum. In general, continuous variables

64    present more statistical power as compared to ordinal or categorical data, and this is likely

65    the case with VAS as compared to ORS. The downside of the conventional VAS is the

66    difficulty to train raters to score different gradations consistently, and as observed by de

67    Grauw and van Loon (2016), the inter-rater reliability may be low. In this case, the tagged

68    VAS (tVAS), which is a VAS with internal anchors, has been investigated as a tool to

69    combine the advantages of both ORS and VAS (Nalon et al., 2014; Tuyttens et al., 2009).

70    The tags add information to guide raters through different gradations thereby increasing

71    reliability and facilitating the training of raters (Tuyttens et al., 2009).

72    Previous studies assumed equidistant tags to VAS to assess specific indicators of

73    animal welfare based on existing categories used in ORS (Meeremans et al., 2017; Nalon

74    et al., 2014; Rufener et al., 2018; Tuyttens et al., 2009). However, Vieira et al. (2015)

75    challenged this rationale by presenting a non-equidistant characteristic of tags in VAS as

76    with lameness in dairy goats. In this case, tags that are based on existing categories from

77    ORS are expected to be checked with respect to what their correct positions are on the VAS

78    and whether these are spaced equidistantly or not. As with lameness, many other relevant

3

79 welfare problems vary continuously and could be assessed by a continuous scale rather

80 than an ORS. For broiler chickens, contact dermatitis and related measures are considered

81 important animal welfare indicators. They have been systematically scored using ORS in a

82 variety of scoring scales: contact dermatitis (Allain et al., 2009; de Jong et al., 2014; Ekstrand

83 et al., 1998; Haslam et al., 2007; Martland, 1985; Souza et al., 2018; Welfare Quality®,

84 2009) and bird soiling (Dawkins et al., 2004; Elwinger, 1995; Weeks et al., 1994; Welfare

85 Quality®, 2009; Wilkins et al., 2003), for example. Potential improvement in the use of VAS

86 to assess these indicators seems to warrant further studies, especially testing for reliability.

87 Recent studies have compared ORS and VAS, including tVAS, in animal welfare

88 assessment. For example, Vogt et al. (2017) considered VAS reliable to assess the

89 temperament of animals, and both VAS and ORS were considered reliable scales to assess

90 lameness in dairy cattle (Flower and Weary, 2006). Considering the use of tags in VAS,

91 tVAS and 5-point ORS presented similarly high interobserver reliability for the assessment

92 of lameness in sows, but both were better than for 2-point ORS (Nalon et al., 2014). However

93 interobserver reliability was better for the tVAS than for the ORS (Tuyttens et al., 2009) when

94 assessing lameness in dairy cattle. In contrast, Meeremans et al. (2017) observed that use

95 of tVAS did not improve the reliability of the assessment of fish vitality as compared to

96 categorical scoring.

97 Regarding the decision on the best type of scale, the determinant seems to rely on

98 how observers are able to discriminate between the levels of the indicator (Engel et al.,

99 2003). Based on this, we aimed to test the application of ORS and VAS for four broiler

100 chicken welfare indicators. The indicators were contact dermatitis on the breast and

101 abdominal areas (CD), footpad dermatitis (FP), hock burns (HB) and bird soiling (BS). We

102 studied inter-rater reliability, the correlation between the VAS and ORS and amongst the

103 welfare indicators measured with VAS and ORS. We also tested the equidistance of ORS

104 categories in relation to values measured using the VAS.

4

105

## 2. **Material and Methods**

2.1 Ethical statement

This project was approved by the Animal Use Ethics Committee of the Agricultural Campus (n. 079/2015; November 12th, 2015) of the Federal University of Paraná.

110

2.2 Animals, housing and data collection

A total of 1,303 broiler chickens, randomly selected from 10 flocks, was assessed in the State of Paraná, Southern Brazil, from January 9th to 13th 2017. The sampling size of 1,300 birds was calculated considering a maximum error of 5% and 95% confidence interval. The sample was not selected to be representative of bird welfare in Brazilian industrial broiler chicken units. The poultry barns had sidewalls with wire mesh covered by blackout curtains working as dark house (n = 1) or covered by yellow curtains, with natural lighting (n = 9). The farms were selected as a convenience sample according to our objective, which was to test the ordinal and analogue scales. All units had automatic feeders, nipple drinkers, sprinklers, exhaust fans and wood shaving litter, and nine units maintained evaporative cooling systems. Indoor mean temperature in the units at time of the visit was 27.7 ± 1.4 °C. Average broiler house area was 1,540 ± 187 $m^2$ and the number of birds per house was 18,904 ± 2,604, with a stocking density of 36.4 ± 0.9 $kg/m^2$. Birds were male and female Cobb 500®, assessed at 41.3 ± 2.0 days of age. The raters were one animal scientist and two veterinarians, one of them experienced in auditing poultry farms. The non-experienced raters underwent a 4 h classroom instruction about the indicators via picture observation, followed by a 4 h training session at the Federal University of Paraná farm. Scales used on the training sessions were obtained from Souza et al. (2018) and Welfare Quality® (2009). One month after the training, the non-experienced raters were asked to score 13 pictures for FP and 15 pictures for CD and BS to check concordance among them and solve any

5

131 doubts before the experiment. Kendall's coefficient of concordance corrected for ties among

132 raters were 0.89 (P=0.002), 0.79 (P=0.004) and 0.93 (P=0.001) for FP, CD and BS,

133 respectively, and were considered adequate (Landis and Koch, 1977).

134 Raters scored each bird simultaneously but independently. They performed a visual

135 inspection of a total 130 birds from five locations in each poultry house. The feet of the birds

136 were cleaned by gently rubbing with the tip of observer's fingers. All assessors scored each

137 bird simultaneously, so that any lesions were seen by all immediately after the cleaning

138 procedure. Following the regular Welfare Quality procedures, birds were not individually

139 identified. As for CD and BS, the original ORS by Souza et al. (2018) were applied, which

140 included a colour picture and a description of each level of the scale. For FP and HB, the

141 scales by the Welfare Quality® (2009) were used, including a colour picture representative

142 of each level of the scale (Fig. 1). To collect data, a questionnaire was developed at the

143 QuickTapSurvey® website to be used as a mobile phone application. Raters scored each

144 bird using both the ORS and the VAS for each indicator. The application presented the ORS

145 followed by VAS, thus the raters usually scored ORS first.. In the ORS, the raters had to

146 select a score on a 4- or 5-point scale. The VAS consisted of a line initially designed with 10

147 cm, and proportional to this length depending on the screen size, and anchored only with

148 the minimum and the maximum score at each end (absence or severe CD, FP, HB and BS).

149 The raters could move a marker along the line to register the level of severity observed in

150 the bird for each indicator. Data from QuickTapSurvey® were downloaded into an Excel file

151 and checked for errors before use.

152

153 [Insert Fig. 1]

154 Fig. 1. Ordinal scales for the assessment of four broiler chicken welfare indicators; [1] (Souza
155 et al., 2018), [2] (Welfare Quality®, 2009).

156

6

157    2.3 Statistical analysis

158         Linear mixed models were fitted to evaluate the inter-rater reliability. The intraclass

159    correlation coefficient (ICC) was estimated as a measure of inter-rater reliability, and

160    bootstrap confidence intervals were derived based on 5,000 simulations. The total data

161    variability (TDV) was decomposed into variability attributed or not attributed to the raters

162    (VNA). ICC values were calculated based on the VNA:TDV ratio, adjusted for the variability

163    between poultry farms.Poultry farm was included as a factor in the model, as measurements

164    in individual animals within the same farm tended to be related to each other. Thus, for any

165    animal welfare indicator, indicated as Y, the following linear mixed model was defined:

166

167                         $$Y_{ijkl} = n + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \gamma_k + \epsilon_{ijkl}$$ , where

168

169    $Y_{ijkl}$ is the *l-th* assessment of the rater *j* in the animal *k* of the poultry farm *i;*

170    $\alpha_i \sim Normal(0, \sigma_\alpha^2)$ is the random effect of poultry farm;

171    $\beta_j \sim Normal(0, \sigma_\beta^2)$ is the random effect of rater;

172    $\gamma_k \sim Normal(0, \sigma_\gamma^2)$ is the random effect of animal;

173    $(\alpha\beta)_{ij} \sim Normal(0, \sigma_{\alpha\beta}^2)$ is the random effect of the interaction between rater and poultry

174    house;

175    η is the model intercept;

176    $\epsilon_{ijkl} \sim Normal(0, \sigma^2)$ is the random error.

177         Based on this, the ICC was calculated as:

178

179                         $$ICC = \frac{\sigma_\alpha^2 + \sigma_\gamma^2}{\sigma_\alpha^2 + \sigma_\beta^2 + \sigma_{\alpha\beta}^2 + \sigma_\gamma^2 + \sigma^2}$$

180

181    ICC values were estimated under both scales (ORS and VAS), and the difference

182    $ICC_{ORS} - ICC_{VAS}$ was calculated. To evaluate the statistical significance of this difference,

183    the null ==hypothesis== of equality was tested through additional simulations, and the simulated

184    p-values are presented.

185    Spearman's rank correlation coefficients for the mean of values given by the three

186    raters and for the individual values of each rater was used to test correlations between ORS

187    and VAS for all indicators, as well as correlations amongst all indicators measured using the

188    ORS. Pearson Correlation Coefficient was used to test correlations amongst all indicators

189    measured using the VAS. Correlations from 0.3 to 0.6 were considered moderate, and

190    values above 0.6 were considered high (de Jong et al., 2015).

191    Linear mixed models were also fitted to test the assumption of equidistance of ORS

192    categories according to values measured using the VAS. For each indicator, the VAS values

193    were considered as the response variable, and the ORS values as the predictor

194    (independent variables). Random effects of animal, rater, poultry house and interaction

195    between rater and poultry house were also included in the ==models==. Two linear mixed models

196    were fitted for each indicator, assuming (Model 1) or not assuming equidistance (Model 2)

197    between the scores. In Model 1, ORS was included as a numerical variable defined by the

198    p+1 different values. In Model 2, ORS was included as a categorical variable, not assuming

199    a fixed increment across scores. So, the following models were considered:

200    Model 1: $Y_{ijkl} = n + NRS_{ijkl} + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \gamma_k + \epsilon_{ijkl}$, where $Y_{ijkl}$ and $NRS_{ijkl}$

201    correspond to the rater *j* in the animal *k* of the poultry farm *i* for the *l-th* time in the scales

202    VAS and ORS, respectively;

203    Model                                                                                2:

$Y_{ijkl} = n + \tau_1 \times I(NRS_{ijkl} = 1) + \tau_2 \times I(NRS_{ijkl} = 2) + \tau_3 \times I(NRS_{ijkl} = 3) + \alpha_i + \beta_j + $
204    $(\alpha\beta)_{ij} + \gamma_k + \epsilon_{ijkl}$

8

205    , where $I(NRS_{ijkl} = x)$ is the indicator function, assuming value zero when ORS score is

206    different of an $x$ value, and assuming value one when ORS score is equal to $x$; $\tau_1$, $\tau_2$ and $\tau_3$

207    are the effects that reflect the association between ORS and VAS.

208            To evaluate the equidistance hypothesis, the fitted linear mixed models were

209    compared using the Akaike Information Criterion (AIC), following the method of Burnham

210    and Anderson (2002), based on the evidence ratio, defined by:

211

$$ER = \frac{1}{\exp\left(-0.5(AIC_{mod1} - AIC_{mod2})\right)}$$

212

213

214     ER will be equal to 1 if the evidence for both models is the same. The greater ER, the

215    greater the evidence for model 2 (non-equidistant), whereas ER moves toward zero if model

216    1 (equidistance) has the highest evidence.

217            For each indicator for which the ER analysis confirmed non-equidistance, the best

218    values for the non-equidistant tags of VAS were obtained through classification tree analysis,

219    with VAS as the predictor and ORS the response. The classification tree method proposed

220    by Breiman et al. (1984) employs successive partitions of a sample to constitute subsamples

221    that are homogeneous in relation to response values, in our case ORS. The rules for the

222    partitions were as $VAS < x$ versus $VAS \geq x$, being a VAS value, so that the observations

223    were allocated to different subsamples (nodes) according to the partition rules. The final

224    number of nodes was defined based on a cross validated procedure. In addition, the number

225    of tags for the ordinal scale was also considered to determine the number of final nodes in

226    a proper way. Analyses were performed using R Statistical Computing Environment

227    software version 3.3.1 (R Core Team, 2016), through the packages boot (Angelo Canty and

228    Brian Ripley, 2016), lme4 (Bates et al., 2015) and rpart (Therneau et al., 2015).

9

229

230    3.  **Results**

231        Estimated inter-rater reliability was higher for CD and HB using VAS, and higher for

232    FP using ORS (Table 1).

233    Table 1. Estimates of inter-rater reliability and confidence interval, 5,000 bootstrap samples,

234    for animal welfare indicators from 1,303 broiler chickens assessed on farm by three raters

235    using both ordinal scale (ORS) and visual analogue scale (VAS).

| Welfare indicator | Scale | Intraclass correlation | Confidence interval (95%) | P-value* |
|---|---|---|---|---|
| Contact dermatitis on the breast and abdominal areas | ORS | 0.68 | (0.58 - 0.77) | <0.001 |
| | VAS | 0.77 | (0.67 - 0.85) | |
| Footpad dermatitis | ORS | 0.91 | (0.87 - 0.93) | <0.001 |
| | VAS | 0.88 | (0.83 - 0.92) | |
| Hock burns | ORS | 0.67 | (0.55 - 0.76) | <0.001 |
| | VAS | 0.72 | (0.60 - 0.80) | |
| Bird soiling | ORS | 0.61 | (0.46 - 0.73) | 0.447 |
| | VAS | 0.54 | (0.36 - 0.69) | |

236    *ORS x VAS intraclass correlation

237        High correlations were observed between ORS and VAS for each welfare indicator,

238    considering mean and individual values (Table 2). When indicators were correlated amongst

239    them, within each scale, we observed similar level of correlation of data using ORS and VAS

240    (Table 2).

241

242    Table 2a. Correlation of ordinal scale (ORS) and visual analogue scale (VAS) for the mean

243    of values given by the three raters and for the individual values of each rater.

244    Table 2b. Correlation of broiler chicken welfare indicators measured on farm using ORS and

245    VAS, 1,303 birds.

| Table 2a | | | Table 2b | | | | |
|---|---|---|---|---|---|---|---|
| Indicator | Spearman rank correlation between ORS and VAS* | | Correlation between indicators* (ORS, Spearman correlation; VAS, Pearson correlation) | | | | |
| | Mean | Individual | Indicator | Scale | FP | HB | BS |

10

| Indicator | | | | | | | |
|---|---|---|---|---|---|---|---|
| Contact dermatitis on the breast and abdominal areas (CD) | 0.96 | 0.89 | CD | ORS VAS | 0.06 0.09 | 0.24 0.35 | 0.34 0.34 |
| Footpad dermatitis (FP) | 0.97 | 0.95 | FP | ORS VAS | | 0.17 0.26 | 0.08 0.12 |
| Hock burn (HB) | 0.90 | 0.77 | HB | ORS VAS | | | 0.25 0.24 |
| Bird soiling (BS) | 0.94 | 0.81 | | | | | |

246 *$P < 0.0001$

247

248    For all indicators, the strength of evidence for the Model 2, which does not assume

249 equidistance between tags, was higher than 0.99. Thus, the tags on VAS that better

250 represent ORS are not evenly spaced. The calculated tags for each indicator are shown in

251 Fig. 2. The prevalence of absence of soiling (score 0) among the broiler chickens assessed

252 in our study was 0.1 %, while severe HB and FP (score 4) was observed in 1.0% and 4.2%,

253 respectively. Since these frequencies did not allow an adequate tag calculation, scores 0

254 and 1 were aggregated for BS, as well as scores 3 and 4 for HB and FP (Fig. 2).

255

256 [Insert Fig. 2}

257 Fig. 2. Tags for ordinal scale (ORS) for broiler chicken welfare indicators calculated by the

258 classification tree considering visual analogue scale (VAS) as predictor. Percentages refer

259 to the number of birds classified in each ORS category, data from 1,303 birds assessed on

260 farm by three raters.

261  4. **Discussion**

262    Higher ICC for CD and HB using VAS, and for FP using ORS were observed;

263 however, common ICC interpretation suggest that both scales were reliable to assess the

264 animal-based indicators proposed in this study. This warrants further research comparing a

265 greater number of raters. Direct comparison across studies using ORS and VAS is not

266 possible due to different methods employed to estimate reliabilities. For those studies in

11

267 which the reliability was given by a value between 0 and 1, the range of reported values

268 considered to be reliable was similar to the range observed in our study (Flower and Weary,

269 2006; Meeremans et al., 2017; Nalon et al., 2014). As a general guideline, ICC reliability as

270 measured with ICC is considered good when between 0.60 and 0.74, and excellent when

271 higher than 0.75 (Cicchetti, 1994). In the case of BS, lack of difference between ORS and

272 VAS seems related to high data variability. FP is observed as a clearer indicator, perhaps

273 as consequence of a simpler scale. In the case of CD and BS, pictures needed to expose

274 other animal parts, like skin, foot, and feathers, which may induce raters to reflect more

275 about animal condition. In this case, data obtained may be influenced by something else,

276 like experience or personal views (Meagher, 2009).

277     Other factors may have affected inter-rater reliability, such as place of assessment,

278 training, quality of the descriptive textual and photographic material to support the

279 assessment, and the limited number of raters. Studies comparing ORS and VAS for animal

280 welfare purposes frequently combine video recordings and a large group of raters (e.g.

281 Tuyttens et al., 2009; Nalon et al., 2014). In our study, on-farm assessments may have

282 improved inter-rater reliability, even with three raters, since they could have chosen the best

283 angle and touched the birds during the physical assessment. Touching the birds was

284 important to remove dirt to confirm the presence and size of FP and HB. Since only one

285 rater was experienced in broiler chicken welfare assessment, training, rather than

286 experience, may have played an important role in helping raters to discriminate between the

287 levels of each indicator (Meeremans et al., 2017). In addition, successful learning depends

288 on a scoring system with clear definitions and photographs (Gibbons et al., 2012). In our

289 case, training was done with the available scientifically validated scales to score the four

290 proposed indicators. These materials were related to the use of ORS, which means that

291 raters were trained to recognize four or five different levels of severity, depending on the

292 indicator. Nevertheless, raters were able to coherently score birds using the VAS. The

12

quality of the scoring system is important to provide all information required by the raters before and during the assessment, and clear definitions are essential to make scoring systems less dependent on personal experience or any factor that reduces inter-rater reliability (Meagher, 2009). In this regard, it is expected that more comprehensive training material, with pictures of various gradations in severity along the VAS, will increase inter-rater reliability.

Indicators showed the same level of correlation between them, regardless of the type of scale. The exception was the correlation between CD and HB, which was slightly higher when using the VAS compared with the ORS. Both CD and HB had higher inter-rater reliability using VAS, thus probably there was a refinement of the scoring using VAS, which impacted on the correlation between CD and HB. We expected higher correlation between CD, FP and HB, since contact dermatitis has been reported as to develop sequentially on different part of the body, starting with feet and followed by hocks and breast, as bird activity decreases (de Jong et al., 2014; Greene et al., 1985). Other factors, such as early age of modern fast-growing broiler chickens at slaughter and litter quality, may challenge the correlation between different types of skin lesion (Souza et al., 2018). Despite low to moderate correlation between indicators, the type of scale did not affect data interpretation for the selected outcomes in this study, suggesting both scales could be used to assess birds.

High correlation between ORS and VAS for all indicators may suggest applicability of both scales and is in line with results of comparisons between ORS and VAS for pain assessment (Hjermstad et al., 2011). Similar to Flower and Weary (2006), raters were able to coherently transpose ordinal scores into continuous scores even in the absence of internal tags on VAS. One possible limitation of this study was the application of both scales concomitantly, which may have motivated raters to virtually divide the VAS according to the ORS. Equidistant data would support this rationale, as observed by Engel et al. (2003).

13

319 However, data obtained in our study were not equidistant. The lack of equidistance has been

320 observed in other studies using VAS to assess lameness (Thomsen et al., 2008; Vieira et

321 al., 2015; Welsh et al., 1993) and in a study to determine cut-off points in a VAS for pain in

322 patients with chronic musculoskeletal pain (Boonstra et al., 2014). Our results show that the

323 decision regarding the location of tags had direct implication on the number of animals

324 classified in each level of severity. As example, some birds who were scored as 0 using

325 ORS, meaning absence of CD, FP and HB, received grades up to 16 or 20 mm using the

326 VAS. These results are probably indicating that birds had less severe lesions than the ones

327 described on level 1 of the ORS, and the rater had to choose between 0 and 1. In this case,

328 the VAS was more sensitive to allow the rater to choose the best position between 0 and 1.

329 In this example, the number of birds considered clinically absent of CD, FP and HB differed

330 between ORS and VAS.  If the different values of a specific welfare indicator are biased to

331 one extreme, either more concentrated on the higher severity or the lower severity end,

332 scales presenting more detailed assessment, as VAS, may offer higher accuracy.

333 According to Averbuch and Katzper (2004) and Nalon et al. (2014) inserting internal

334 tags on VAS allows combining characteristics of the ORS, and improves uniformity of

335 interpretation, with the flexibility of VAS to identify small changes between the tags. Although

336 the VAS had a high reliability in this study, it is expected that the internal anchors of a tVAS

337 will enable raters to score even more reliably. The position of the internal tags in a tVAS is

338 important because it affects the number of animals in each level. As observed in Fig. 2,

339 categories 0 and 1 were often narrower than the more severe categories. Perhaps the ORS

340 over-emphasizes the milder cases, which were the most common for three indicators, while

341 the VAS allows raters to better differentiate between the scores. Thus, to compare ORS and

342 tVAS, it is important to have clear definitions about the position of different ORS categories

343 along the continuous scale, and raters should be clearly instructed and trained on how to

344 use the scale. This issue deserves more attention and seems especially relevant depending

14

345    on the goal of the assessment, which may be to provide best practice recommendations or

346    may be associated with sanctions (Main and Mullan, 2012) or bonuses for certification

347    processes.

348       Many studies have been done to encourage the adoption of regular broiler chicken

349    welfare assessment worldwide. This permanent monitoring of welfare may include the use

350    of correlations, such as of contact dermatitis on farm and at the slaughterhouse (de Jong et

351    al., 2015), as well as the use of technology to automate assessment on farm or at the

352    slaughterhouse (Sassi et al., 2016). FP has been accepted as an important welfare indicator

353    for surveillance purposes (European Commission, 2017), and automation of this

354    assessment seems a priority. For automated assessment through image analysis, the ORS

355    are commonly used, and in the case of FP they seem adequate. When both VAS and ORS

356    work well, the choice of the scale will include a critical analysis of the conditions related to

357    their use (Hjermstad et al., 2011). Adoption of an animal welfare indicator by organizations

358    will depend on reliability, validity, sensitivity and power, but also feasibility and efficiency.

359    VAS, including tVAS, presents potential to be considered for different animal welfare

360    strategies, in addition to animal welfare assessment. As example, it may be used to validate

361    automated monitoring of indicators showing higher inter-rater reliability using VAS or, since

362    VAS is more sensitive (Welsh et al., 1993). Application may include its use during

363    inspections for certification processes and as part of a verification procedure in an animal

364    welfare management system (Souza and Molento, 2018), in studies in which high sensitivity

365    is needed; or tVAS may be used as a silver standard for automated monitoring tools, since

366    it is more likely to detect small differences and changes along time. In addition, future work

367    studying the biological validity of both VAS and ORS with appropriate standards, as for

368    instance histological assessments to check dermatitis severity, seem warranted to further

369    understand accuracy of the measurements.

370

15

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

371    5.  **Conclusion**

372        This is the first study to compare ORS and VAS for the selected broiler chicken

373    welfare indicators. Both ORS and VAS were considered reliable to assess the broiler

374    chicken welfare indicators CD, FP, HB and BS, despite some differences in inter-rater

375    reliability. Although higher inter-rater reliability may lead to refined correlation studies, the

376    interpretation of correlation did not differ between VAS and ORS. VAS, including tVAS,

377    presents potential to add sensitivity on animal welfare assessment, and is a tool to be further

378    explored in validation and certification protocols, especially in studies in which high

379    sensitivity is needed. In this case, considering that results from animal welfare assessment

380    may have direct implications to the animals and other stakeholders, the use of tVAS will

381    demand clear specification about the position of tags on the continuous scale as well as the

382    training of raters.

383

384    **Acknowledgements**

385

386

387

388

389    **References**

390    Allain, V., Mirabito, L., Arnould, C., Colas, M., Le Bouquin, S., Lupo, C., Michel, V., 2009.

391        Skin lesions in broiler chickens measured at the slaughterhouse: relationships between

392        lesions and between their prevalence and rearing factors. Br. Poult. Sci. 50, 407–417.

393        https://doi.org/10.1080/00071660903110901

394    Angelo Canty, Brian Ripley, 2016. boot: Bootstrap R (S-Plus) Functions.

395    AssureWel, 2014. AssureWel - Advancing Animal Welfare Assurance [WWW Document].

396        Broilers. URL http://www.assurewel.org/broilers (accessed 9.4.18).

16

397   Averbuch, M., Katzper, M., 2004. Assessment of Visual Analog versus Categorical Scale for

398       Measurement of Osteoarthritis Pain. J. Clin. Pharmacol. 44, 368–372.

399       https://doi.org/10.1177/0091270004263995

400   Bates, D., Mächler, M., Bolker, B., Walker, S., 2015. Fitting linear mixed-effects models

401       using lme4. J. Stat. Softw. 67, 1–48. https://doi.org/10.18637/jss.v067.i01

402   Boonstra, A.M., Preuper, H.R.S., Balk, G.A., Stewart, R.E., 2014. Cut-off points for mild,

403       moderate, and severe pain on the visual analogue scale for pain in patients with chronic

404       musculoskeletal              pain.           Pain              155,              2545–2550.

405       https://doi.org/10.1016/j.pain.2014.09.014

406   Breiman, L., Friedman, J., Stone, C.J., Olshen, R.A., 1984. Classification and regression

407       trees. CRC Press, Boca Raton.

408   Burnham, K.P., Anderson, D.R., 2002. Model selection and multimodel inference: a practical

409       information-theoretic approach, 2nd ed. Springer-Verlag, New York.

410   Cicchetti, D. V., 1994. Guidlines, Criteria, and Rules of Thumb for Evalauting Normed and

411       Standardized Assessment Instruments in Psychology. Psychol. Assess. 6, 284–290.

412       https://doi.org/10.1037/1040-3590.6.4.284

413   Dawkins, M.S., Donnely, . A. E., Jones, T.A., 2004. Chicken welfare is influenced more by

414       housing conditions than by stocking density. Nature 427, 342–343.

415   de Grauw, J.C., van Loon, J.P.A.M., 2016. Systematic pain assessment in horses. Vet. J.

416       209, 14–22. https://doi.org/10.1016/j.tvjl.2015.07.030

417   de Jong, I.C., Gunnink, H., van Harn, J., 2014. Wet litter not only induces footpad dermatitis

418       but also reduces overall welfare, technical performance, and carcass yield in broiler

419       chickens. J. Appl. Poult. Res. 23, 51–58. https://doi.org/10.3382/japr.2013-00803

420   de Jong, I.C., Hindle, V.A., Butterworth, A., Engel, B., Ferrari, P., Gunnink, H., Perez Moya,

421       T., Tuyttens, F.A.M., van Reenen, C.G., 2015. Simplifying the Welfare Quality®

422       assessment protocol for broiler chicken welfare. Animal 10, 117–27.

17

423        https://doi.org/10.1017/S1751731115001706

424    Ekstrand, C., Carpenter, T.E., Andersson, I., Algers, B., 1998. Prevalence and control of

425        foot-pad dermatitis in broilers in Sweden. Br. Poult. Sci. 39, 318–24.

426        https://doi.org/10.1080/00071669888845

427    Elwinger, K., 1995. Broiler production under varying population densities - A field study .

428        Arch. fur Geflugelkd. 59, 209–215.

429    Engel, B., Bruin, G., Andre, G., Buist, W., 2003. Assessment of observer performance in a

430        subjective scoring system: Visual classification of the gait of cows. J. Agric. Sci. 140,

431        317–333. https://doi.org/10.1017/S0021859603002983

432    European Commission, 2017. Study on the application of the broilers directive (DIR

433        2007/43/EC) and development of welfare indicators. Brussels.

434        https://doi.org/10.1149/1.3477934

435    Fleming, P.A., Clarke, T., Wickham, S.L., Stockman, C.A., Barnes, A.L., Collins, T., Miller,

436        D.W., 2016. The contribution of qualitative behavioural assessment to appraisal of

437        livestock welfare. Anim. Prod. Sci. 56, 1569–1578. https://doi.org/10.1071/AN15101

438    Flower, F.C., Weary, D.M., 2006. Effect of Hoof Pathologies on Subjective Assessments of

439        Dairy Cow Gait. J. Dairy Sci. 89, 139–146. https://doi.org/10.3168/jds.S0022-

440        0302(06)72077-X

441    Gibbons, J., Vasseur, E., Rushen, J., De Passillé, A.M., 2012. A training programme to

442        ensure high repeatability of injury scoring of dairy cows. Anim. Welf. 21, 379–388.

443        https://doi.org/10.7120/09627286.21.3.379

444    Global Animal Partnership's, 2018. Animal Welfare Rating Standard For Chickens Raised

445        for Meat v3.1 [WWW Document]. v3.1. URL https://globalanimalpartnership.org/wp-

446        content/uploads/2018/04/GAP-Standard-for-Meat-Chickens-v3.1-20180403.pdf

447        (accessed 9.4.18).

448    Greene, J.A., McCracken, R.M., Evans, R.T., 1985. A contact dermatitis of broilers -clinical

18

449      and      pathological      findings.      Avian      Pathol.      14,      23–38.

450      https://doi.org/10.1080/03079458508436205

451 Grosso, L., Battini, M., Wemelsfelder, F., Barbieri, S., Minero, M., Dalla Costa, E., Mattiello,

452      S., 2016. On-farm Qualitative Behaviour Assessment of dairy goats in different housing

453      conditions.      Appl.      Anim.      Behav.      Sci.      180,      51–57.

454      https://doi.org/10.1016/j.applanim.2016.04.013

455 Haslam, S.M., Knowles, T.G., Brown, S.N., Wilkins, L.J., Kestin, S.C., Warriss, P.D., Nicol,

456      C.J., 2007. Factors affecting the prevalence of foot pad dermatitis, hock burn and

457      breast burn in broiler chicken. Br. Poult. Sci. 48, 264–275.

458 Hjermstad, M.J., Fayers, P.M., Haugen, D.F., Caraceni, A., Hanks, G.W., Loge, J.H.,

459      Fainsinger, R., Aass, N., Kaasa, S., 2011. Studies comparing numerical rating scales,

460      verbal rating scales, and visual analogue scales for assessment of pain intensity in

461      adults: A systematic literature review. J. Pain Symptom Manage. 41, 1073–1093.

462      https://doi.org/10.1016/j.jpainsymman.2010.08.016

463 Landis, J.R., Koch, G.G., 1977. The measurement of observer agreement for categorical

464      data. Biometrics 33, 159–174. https://doi.org/10.2307/2529310

465 Main, D.C.J., Mullan, S., 2012. Economic, education, encouragement and enforcement

466      influences      within      farm      assurance      schemes.      Anim.      Welf.      21,      107–111.

467      https://doi.org/10.7120/096272812X13345905673881

468 Martland, M.F., 1985. Ulcerative dermatitis dm broiler chickens: the effects of wet litter. Avian

469      Pathol. 14, 353–364. https://doi.org/10.1080/03079458508436237

470 Meagher, R.K., 2009. Observer ratings: Validity and value as a tool for animal welfare

471      research.      Appl.      Anim.      Behav.      Sci.      119,      1–14.

472      https://doi.org/10.1016/j.applanim.2009.02.026

473 Meeremans, P., Yochum, N., Kochzius, M., Ampe, B., Tuyttens, F.A.M., Uhlmann, S.S.,

474      2017. Inter-rater reliability of categorical versus continuous scoring of fish vitality: Does

19

475    it affect the utility of the reflex action mortality predictor (RAMP) approach? PLoS One

476    12, 1–22. https://doi.org/10.1371/journal.pone.0179092

477  Minero, M., Dalla, E., Dai, F., Anne, L., Murray, M., Canali, E., Wemelsfelder, F., 2016. Use

478    of Qualitative Behaviour Assessment as an indicator of welfare in donkeys. Appl. Anim.

479    Behav. Sci. 174, 147–153.

480  Nalon, E., Maes, D., Van Dongen, S., van Riet, M.M.J., Janssens, G.P.J., Millet, S.,

481    Tuyttens, F.A.M., 2014. Comparison of the inter- and intra-observer repeatability of

482    three      gait-scoring      scales      for      sows.      Animal      8,      650–659.

483    https://doi.org/10.1017/S1751731113002462

484  R Core Team, 2016. A language and environment for statistical computing.

485  Rufener, C., Baur, S., Stratmann, A., Toscano, M.J., 2018. A Reliable Method to Assess

486    Keel Bone Fractures in Laying Hens From Radiographs Using a Tagged Visual

487    Analogue Scale. Front. Vet. Sci. 5, 1–8. https://doi.org/10.3389/fvets.2018.00124

488  Sassi, N. Ben, Averós, X., Estevez, I., 2016. Technology and poultry welfare. Animals 6, 1–

489    21. https://doi.org/10.3390/ani6100062

490  Souza, A., Soriano, V., Schnaider, M., Rucinque, D., Molento, C., 2018. Development and

491    refinement of three animal-based broiler chicken welfare indicators. Anim. Welf. 27,

492    263–274.

493  Souza, A.P.O., Molento, C.F.M., 2018. Proposal of a management system to develop an

494    animal   welfare   strategy   for   the   animal   food   chain.   CAB   Rev.   13,   1–11.

495    https://doi.org/10.1079/PAVSNNR201813001

496  Therneau, T., Atkinson, B., Brian Ripley, 2015. rpart: Recursive Partitioning and Regression

497    Trees.

498  Thomsen, P.T., Munksgaard, L., Tøgersen, F.A., 2008. Evaluation of a Lameness Scoring

499    System for Dairy Cows. J. Dairy Sci. 91, 119–126. https://doi.org/10.3168/jds.2007-

500    0496

20

501    Tuyttens, F.A.M., de Graaf, S., Heerkens, J.L.T., Jacobs, L., Nalon, E., Ott, S., Stadig, L.,

502    Van Laer, E., Ampe, B., 2014. Observer bias in animal behaviour research: Can we

503    believe what we score, if we score what we believe? Anim. Behav. 90, 273–280.

504    https://doi.org/10.1016/j.anbehav.2014.02.007

505    Tuyttens, F.A.M., Sprenger, M., Van Nuffel, A., Maertens, W., Van Dongen, S., 2009.

506    Reliability of categorical versus continuous scoring of welfare indicators: Lameness in

507    cows as a case study. Anim. Welf. 18, 399–405.

508    https://doi.org/10.1016/j.applanim.2010.05.003

509    Vieira, A., Oliveira, M.D., Nunes, T., Stilwell, G., 2015. Making the case for developing

510    alternative lameness scoring systems for dairy goats. Appl. Anim. Behav. Sci. 171, 94–

511    100. https://doi.org/10.1016/j.applanim.2015.08.015

512    Vogt, A., Aditia, E.L., Schlechter, I., Schütze, S., Geburt, K., Gauly, M., König von Borstel,

513    U., 2017. Inter- and intra-observer reliability of different methods for recording

514    temperament in beef and dairy calves. Appl. Anim. Behav. Sci. 195, 15–23.

515    https://doi.org/10.1016/j.applanim.2017.06.008

516    Weeks, C.A., Nicol, C.J., Sherwin, C.M., Kestin, S.C., 1994. Comparison of the behaviour

517    of broiler chickens in indoor and free-range environments. Anim. Welf. 3, 179–192.

518    Welfare Quality®, 2009. Welfare Quality ® Assessment protocol for poultry (broilers, laying

519    hens). Welfare Quality Consortium, Lelystad, The Netherlands, p. 116.

520    Welsh, E.M., Gettinby, G., Nolan, A.M., 1993. Comparison of a visual analogue scale and a

521    numerical rating scale for assessment of lameness, using sheep as a model. Am. J.

522    Vet. Res. 54, 976–983.

523    Wemelsfelder, F., Hunter, T. E., Mendl, M., Lawrence, A. B., 2001. Assessing the 'whole

524    animal': A free choice profiling approach. Anim. Behav. 62, 209–220.

525    Wilkins, L.J., Brown, S.N., Phillips, A.J., Warriss, P.D., 2003. Cleanliness of broilers when

526    they arrive at poultry processing plants. Vet. Rec. 701–703.

21

1
2
3        527
4
5        528
6
7        529
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58                                                                                        22
59
60