

On Jones et al.'s method for extending Bland-Altman plots to limits of agreement with the mean for multiple observers

Heidi S. Christensen^{1,2,3}, Jens Borgbjerg⁴, Lars Børty² and Martin Bøgsted^{1,2,3*}

¹Department of Clinical Medicine, Aalborg University, Aalborg, Denmark

²Department of Haematology, Aalborg University Hospital, Aalborg, Denmark

³Clinical Cancer Research Center, Aalborg University Hospital, Aalborg, Denmark

⁴Department of Radiology, Aarhus University Hospital, Aarhus, Denmark

*Corresponding author, email: martin.boegsted@rn.dk

Abstract

Background To assess the agreement of continuous measurements between a number of observers, Jones et al. introduced limits of agreement with the mean (LOAM) for multiple observers, representing how much an individual observer can deviate from the mean measurement of all observers. Besides the graphical visualisation of LOAM, suggested by Jones et al., it is desirable to supply LOAM with confidence intervals and to extend the method to the case of multiple measurements per observer.

Methods We reformulate LOAM under the assumption the measurements follow an additive two-way random effects model. Assuming this model, we provide estimates and confidence intervals for the proposed LOAM. Further, this approach is easily extended to the case of multiple measurements per observer.

Results The proposed method is applied on two data sets to illustrate its use. Specifically, we consider agreement between measurements regarding tumour size and aortic diameter. For the latter study, three measurement methods are considered.

Conclusions The proposed LOAM and the associated confidence intervals are useful for assessing agreement between continuous measurements.

Keywords Accuracy, limits of agreement with the mean, continuous measurements

27 **1 Background**

28 Clinical decisions regarding diagnosis or treatment are often based on one or more measured quantities
29 such as blood pressure, tumour size, or the diameter of an aorta. To understand the limitations of using such
30 measurements in clinical practice, it is important to quantify how much the measurements may vary.

31 For almost three decades, Bland-Altman plots have been the standard method for graphical assessment
32 of agreement between continuous measurements made by two observers or methods on a number of
33 subjects [1]. In particular, Bland-Altman plots are often used to assess how well a new measurement method
34 compares to a current standard method. However, if the goal is to assess the variability of measurements
35 made by different observers it is preferable to consider more than two observers.

36 This prompted Jones et al. to suggest an extension of Bland-Altman's graphical method for assessing
37 *limits of agreement between two observers* to the *limits of agreement with the mean (LOAM) for multiple*
38 *observers* [2]. Jones et al.'s LOAM have the advantage that they quantify agreement between measurements
39 on the same scale as the measurements themselves, in contrast to the intra-class correlation (ICC) that has
40 no unit of measure and always takes value between 0 and 1.

41 In more detail, consider a study where a continuous quantity is observed on a subjects by b observers
42 (or methods). We let y_{ij} denote an observation from a random variable Y_{ij} , which models the measurement
43 performed on the i^{th} subject by the j^{th} observer for $i = 1, \dots, a$ and $j = 1, \dots, b$. Assuming no preferred
44 observer, Jones et al. suggested to assess the agreement between measurements made by different
45 observers by investigating how much the measurements vary around the subject-specific average [2]. More
46 formally, they were interested in how much the differences $D_{ij} = Y_{ij} - \bar{Y}_i$ are likely to vary, where
47 \bar{Y}_i denotes the average measurement for subject i across the b observers. For visualising the data, Jones et
48 al. propose to consider a plot of the observed differences $d_{ij} = y_{ij} - \bar{y}_i$ against the observed subject-
49 specific average \bar{y}_i . We will refer to this as an *agreement plot*. For an example of an agreement plot see
50 Figure 1 below. An agreement plot can, for example, help to detect whether the spread of the differences is
51 associated to the size of the measurements, or, at least when a and b are not too large, whether some
52 observers tend to always make large, small, or more varying measurements.

53 Further, Jones et al. equipped the agreement plot with horizontal lines representing the estimated
54 95% LOAM, which are given by $\pm 1.96s$, where s is the estimate of the residual standard deviation in a
55 two-way analysis of variance (ANOVA) including subject and observer as fixed effects. Thus, s is only a
56 measure of the residue variation left after accounting for possible subject and observer effects. On one
57 hand, if there is a non-negligible observer effect, this should be included in the variability of the
58 differences d_{ij} when constructing the LOAM. On the other hand, in the (unrealistic) case of no variation
59 due to observer the 95% LOAM lines suggested by Jones et al. are biased and inefficiently estimated, as it
60 would be custom to refit the ANOVA model without the adjustment for observer effect and adjust the
61 degrees of freedom for s accordingly.

62 In conclusion, although the method has gained an increasing interest over the years, Jones et al. did not
63 provide a way to: 1) assess the variation of the LOAM estimate, 2) integrate variation due to different
64 observers, and 3) extend the method to multiple observations per observer.

65 In this paper, we suggest formalising Jones et al.'s approach under a simple two-way random effects
66 model which allows us to formulate a coherent statistical inference procedure for the LOAM. In addition,
67 we provide not only an implementation in the statistical programming software R, but also simple formulae
68 which can be implemented in, e.g., statistical programming languages, Excel, or automatic web-modules for
69 data collection.

70

71 **2 Methods**

72 **2.1 A revised version of the limits of agreement with the mean**

73 We propose to derive LOAM assuming a random effects model for the measurements. Assuming a statistical
74 model provides a theoretical framework in which the LOAM can be constructed in a transparent way and
75 furthermore enables us to supply estimates and confidence intervals (CIs) for the LOAM.

76 **2.1.1 Statistical model**

77 In the following we assume the measurements follow a two-way random effects model given by

$$Y_{ij} = \mu + A_i + B_j + E_{ij}, \quad (1)$$

78 where μ describes the overall mean, and A_i , B_j , and E_{ij} are independent random variables following zero-
 79 mean normal distributions with variances σ_A^2 , σ_B^2 , and σ_E^2 , respectively.

80 Under this model, measurements made by different observers are uncorrelated if they are on different
 81 subjects, while they are positively correlated with covariance σ_A^2 for the same subjects. Further, the
 82 covariance between measurements made by the same observer for different subjects is σ_B^2 . Note that the
 83 measurements are assumed to be homoscedastic, i.e. has common variance, where the common variance is
 84 given by $\sigma_A^2 + \sigma_B^2 + \sigma_E^2$. That is, the variance is split into three components: the inter-subject, inter-observer,
 85 and residual variance. Here we follow the convention of referring to the residual variance σ_E^2 as the intra-
 86 observer variance. Further, note that we assume a balanced data setup, where each observer has evaluated
 87 all the subjects.

88 2.1.2 Proposed limits of agreement with the mean

89 Under the two-way random effects model stated in Eq. (1), the difference between an individual
 90 measurement and the subject-specific mean, D_{ij} , is normally distributed with mean zero and variance $(\sigma_B^2 +$
 91 $\sigma_E^2)(b - 1)/b$. Thus, under this model we expect 95% of these differences to be within the limits

$$\pm 1.96 \sqrt{\frac{b-1}{b} (\sigma_B^2 + \sigma_E^2)}. \quad (2)$$

92 We propose the above as the 95% LOAM.

93 To estimate σ_B^2 and σ_E^2 under the suggested two-way random effects model, we use the unbiased and
 94 consistent ANOVA estimates (see, e.g., Chapter 4 of Searle et al. [3]), given by

$$\hat{\sigma}_B^2 = \frac{MSB - MSE}{a}, \quad \hat{\sigma}_E^2 = MSE, \quad (3)$$

95 where $MSB = SSB/\nu_B$ and $MSE = SSE/\nu_E$, with $SSB = a \sum_{j=1}^b (\bar{y}_{.j} - \bar{y}_{..})^2$ and $SSE = \sum_{i=1}^a \sum_{j=1}^b (y_{ij} - \bar{y}_{i.} -$
 96 $\bar{y}_{.j} + \bar{y}_{..})^2$ denoting the sums of squares for the observer and residual term, and $\nu_B = b - 1$ and $\nu_E = (a -$
 97 $1)(b - 1)$. Further, $\bar{y}_{i.}$, $\bar{y}_{.j}$, and $\bar{y}_{..}$ denote the subject-specific, observer-specific, and overall average,
 98 respectively. Using the estimates of σ_B^2 and σ_E^2 from Eq. (3), we obtain the following estimate of the 95%
 99 LOAM:

$$\pm 1.96 \sqrt{\frac{SSB + SSE}{N}} = \pm 1.96 \sqrt{\frac{\sum_{i=1}^a \sum_{j=1}^b (y_{ij} - \bar{y}_i)^2}{N}}, \quad (4)$$

100 where $N = ab$ is the total number of measurements. For comparison, Jones et al.'s estimate of the LOAM is
 101 given by

$$103 \quad \pm 1.96 \sqrt{\frac{\sum_{i=1}^a \sum_{j=1}^b (y_{ij} - \bar{y}_i - \bar{y}_j - \bar{y})^2}{v_E}} = \pm 1.96 \hat{\sigma}_E,$$

102 which does not include variation due to observers.

104 2.1.3 Confidence intervals

105 Instead of simply reporting the estimated LOAM given by Eq. (4), it is more informative to report CIs.
 106 However, as the distribution of the LOAM is quite complicated, we only supply approximate CIs.

107 Graybill and Wang propose a method for constructing (approximate) efficient CIs for linear
 108 combinations of variances, [4]. To construct CIs for the LOAM in Eq. (2), we first use the method by Graybill
 109 and Wang to construct a CI for the term inside the square root of the LOAM. Next, that CI is transformed into
 110 a CI for the upper LOAM by taking the square root and then multiplying by 1.96 [see Additional file 1 for
 111 details]. The resulting approximate (and asymmetric) 95% CI for the upper 95% LOAM is given by

$$(1.96\sqrt{(SSB + SSE - L)/N}, 1.96\sqrt{(SSB + SSE + H)/N}), \quad (5)$$

112 where

$$113 \quad L = \sqrt{l_B^2 SSB^2 + l_E^2 SSE^2}, \quad H = \sqrt{h_B^2 SSB^2 + h_E^2 SSE^2}$$

114 with $l_x = 1 - 1/F_{0.975; v_x, \infty}$ and $h_x = 1/F_{0.025; v_x, \infty} - 1$ for $x = B$ and $x = E$ (see Graybill and Wang for
 115 other choices of l_x and h_x [4]). Here $F_{\alpha, m, n}$ is the α -quantile for the F -distribution with m numerator and n
 116 denominator degrees of freedom. A 95% CI for the lower 95% LOAM is simply obtained by negation of the
 117 end points of the CI for the upper LOAM, that is,

$$118 \quad (-1.96\sqrt{(SSB + SSE + H)/N}, -1.96\sqrt{(SSB + SSE - L)/N}).$$

119 Simulations under the two-way random effects model in Eq. (1) indicate that the coverage probability
 120 for the approximate CI is in reality quite close to the wanted 95% even with a low number of observers [see
 121 Figure 1 in Additional file 2].

122 2.1.4 Sample size calculations

123 When planning an agreement study, it is often desirable to investigate how many measurements are
 124 necessary to obtain a certain level of precision in terms of a specified width of the CI for the LOAM. From
 125 Eq. (5) it is clear that the value of L and H determine the width of the CI for the LOAM; specifically, the CI
 126 gets narrower as L and H approaches zero. In turn, this happens when b is increased, since l_x and h_x
 127 approaches zero, when v_x increases for both $x = B$ and $x = E$. Thus, to obtain a higher precision we have
 128 to increase the number of observers, b , while it is not enough to increase the number of subjects.

129 Therefore, assume we have a fixed number of subjects a we want to include in a future study to assess
 130 agreement between measurements. To determine the number of observers necessary to obtain a desired
 131 width W of the 95% CI, we require initial estimates of σ_B^2 and σ_E^2 , say $\hat{\sigma}_{B,0}^2$ and $\hat{\sigma}_{E,0}^2$, which can be obtained
 132 from, e.g., a pilot study. Exploiting the relations $SSE = v_E \hat{\sigma}_{E,0}^2$ and $SSB = v_B (a \hat{\sigma}_{B,0}^2 + \hat{\sigma}_{E,0}^2)$, we can express the
 133 width of the CI in Eq. (5) in terms of the variance estimates rather than the sum of squares. Further, we let
 134 the estimates be given by the initial estimates $\hat{\sigma}_{B,0}^2$ and $\hat{\sigma}_{E,0}^2$, and set the width equal to W . That is, we want
 135 to solve the following equation with respect to b :

$$W = \frac{1.96}{\sqrt{N}} \left(\sqrt{v_B (a \hat{\sigma}_{B,0}^2 + \hat{\sigma}_{E,0}^2) + v_E \hat{\sigma}_{E,0}^2 + H_0} - \sqrt{v_B (a \hat{\sigma}_{B,0}^2 + \hat{\sigma}_{E,0}^2) + v_E \hat{\sigma}_{E,0}^2 - L_0} \right), \quad (6)$$

136 where

$$L_0 = \sqrt{l_B^2 v_B^2 (a \hat{\sigma}_{B,0}^2 + \hat{\sigma}_{E,0}^2)^2 + l_E^2 v_E^2 (\hat{\sigma}_{E,0}^2)^2}, \quad H_0 = \sqrt{h_B^2 v_B^2 (a \hat{\sigma}_{B,0}^2 + \hat{\sigma}_{E,0}^2)^2 + h_E^2 v_E^2 (\hat{\sigma}_{E,0}^2)^2}. \quad (7)$$

137 Note that v_B, v_E, l_B, l_E, h_B , and h_E all depend on b . The equation can then be solved numerically with
 138 respect to b to find the number of observers needed to obtain an expected width W of the 95% CI for the
 139 95% LOAM.

140 2.1.5 Inference on the variance components

141 In order to assess the extent of the inter-subject, inter-observer, and intra-observer variations, we suggest
 142 to consider a 95% CI for σ_A, σ_B , and σ_E , respectively.

143 If the ANOVA estimate $\hat{\sigma}_B^2 > 0$, we simply estimate σ_B by $\hat{\sigma}_B = \sqrt{\hat{\sigma}_B^2}$. Using the statistical delta method
 144 [see Additional file 3], we obtain the following approximate 95% CI for σ_B :

$$\hat{\sigma}_B \pm \frac{1.96}{a\hat{\sigma}_B} \sqrt{\frac{(a\hat{\sigma}_B^2 + \hat{\sigma}_E^2)^2}{2\nu_B} + \frac{(\hat{\sigma}_E^2)^2}{2\nu_E}}. \quad (8)$$

145 Results from a small simulation study investigating how well the actual coverage of the approximate
 146 confidence interval matches the desired coverage probability and how this depends on b and the true values
 147 of σ_B and σ_E can be found in the additional files [see Figure 2 in Additional file 2]. In general, the
 148 approximation improves as b increases.

149 It might happen the estimate $\hat{\sigma}_B^2$ is negative due to negative correlation between observations made by
 150 the same observer on different subjects which will indicate a misspecification of the two-way random
 151 effects model formulated in Eq. (1). Negativity can also arise by sampling variation of the unbiased ANOVA
 152 estimates, we have used in this paper. Although it is tempting to suggest setting $\hat{\sigma}_B^2$ to zero in such a case,
 153 this would introduce bias in the estimation. We therefore suggest to report the negative estimates, and
 154 recommend the researcher to comment on the possibility of negatively correlated measurements, and if
 155 that does not seem realistic, to assess whether the CIs are too wide to provide any clinically meaningful
 156 conclusion. It should be assessed whether more observers should be included to improve the precision of
 157 the estimate or whether the model is wrongly specified.

158 As the distribution of $\hat{\sigma}_E^2$ is known in closed form, an exact asymmetric 95% CI can easily be constructed
 159 for σ_E [see Additional file 3] and is given by

$$\left(\hat{\sigma}_E \sqrt{\frac{\nu_E}{\chi_{0.975; \nu_E}^2}}, \hat{\sigma}_E \sqrt{\frac{\nu_E}{\chi_{0.025; \nu_E}^2}} \right), \quad (9)$$

160 where $\hat{\sigma}_E = \sqrt{\hat{\sigma}_E^2}$ and $\chi_{\alpha; \nu_E}^2$ is the α -quantile of a χ^2 -distribution with ν_E degrees of freedom.

161 To provide some context for the scale of $\hat{\sigma}_B$ and $\hat{\sigma}_E$, it may also be constructive to consider $\hat{\sigma}_A = \sqrt{\hat{\sigma}_A^2}$,
 162 where $\hat{\sigma}_A^2 = (MSA - MSE)/b$ is the ANOVA estimate of σ_A^2 where $MSA = SSA/\nu_A$ with $\nu_A = a - 1$ and
 163 $SSA = b \sum_{i=1}^a (\bar{y}_i - \bar{y}.)^2$. The estimate of σ_A may be accompanied by an (approximate) 95% CI, which can
 164 be constructed using the statistical delta method [see Additional file 3]:

$$\hat{\sigma}_A \pm \frac{1.96}{b\hat{\sigma}_A} \sqrt{\frac{(b\hat{\sigma}_A^2 + \hat{\sigma}_E^2)^2}{2\nu_A} + \frac{(\hat{\sigma}_E^2)^2}{2\nu_E}}. \quad (10)$$

165 **2.1.6 Performing an agreement analysis**

166 To investigate agreement between observers, we propose first to make the agreement plot with the
167 estimate and CI for the 95% LOAM from Sections 2.1.2 - 2.1.3, and to calculate the empirical means and
168 standard deviations for the measurements conditional on observer or subject. Inspection of the agreement
169 plot and the empirical means across subject, conditional on observer can be used to reveal whether any
170 observers tend to make unusually large or small measurements. Further, the agreement plot and the
171 conditional empirical standard deviations can be used to check whether the assumption of
172 homoscedasticity of the random model is fulfilled. If the model in Eq. (1) is fitted using statistical software
173 it is often possible to extract residuals and predictions of the observer and subject effects which can be used
174 to check the model assumptions further. Specifically, one may, e.g., consider plots of the residuals against
175 the fitted values, observer number, and subject number, respectively, to further investigate the
176 homoscedasticity assumption. Further, a normal quantile-quantile plot of the residuals as well as of the
177 predictions of the observer and subject effects, respectively, can be used to investigate the normality
178 assumptions. However, if the number of observers or subjects is low, an inspection of how the predictions
179 are distributed may be pointless. See, for example, Section 4.3 in Pinheiro and Bates for a more detailed
180 explanation and illustration of model diagnostics [5]. If it is concluded that the model assumptions are
181 unreasonable, one could consider an appropriate transformation of the data or formulate a variance model
182 to handle heteroscedasticity of the outcome [5] or one could consider using a generalised, linear, and mixed
183 model to handle non-normal distribution of outcomes [6].

184 If the model seems reasonable, we report the estimate and CI for the LOAM. The clinician can then
185 compare the estimated LOAM and associated CI to a clinically acceptable difference between measurements
186 evaluated on the same subject. Whether or not the agreement between measurements is satisfactory
187 depends both on the scale and clinical purpose of the measurements.

188 Next, we may calculate CIs for σ_B and σ_E , and use these along with the point estimates ($\hat{\sigma}_B^2$ and $\hat{\sigma}_E^2$) to
189 compare the order of magnitude of the inter-observer variation with the intra-observer variation. In the
190 rare case where the observer variation is negligible, the observer effect could in principle be removed from
191 the random model, requiring that the CIs for the LOAM are adjusted accordingly [see Additional file 4].

192 The agreement analysis may be supplemented with an estimate and CI for the ICC, which is another
 193 measure for agreement based on the variance components. Various forms of ICCs are listed in McGraw and
 194 Wong for a range of models [5]. The two-way random effects model proposed in this paper corresponds to
 195 Case 2A in McGraw and Wong, with subject as row effect and observer as column effect, and ICC(A, 1) can
 196 then be used to assess absolute agreement of the measurements [5]. The plug-in estimate of ICC(A, 1) is
 197 easily calculated using the estimated variance components:

$$198 \quad ICC(\widehat{A}, 1) = \frac{\hat{\sigma}_A^2}{\hat{\sigma}_A^2 + \hat{\sigma}_B^2 + \hat{\sigma}_E^2}.$$

199 We refer to Table 7 in McGraw and Wong for an approximate CI for ICC(A,1) [5].

200 **2.2 Multiple measurements on each subject per observer**

201 The proposed LOAM and their estimates and CIs can easily be extended to the case where each observer
 202 performs multiple measurements on every subject. If each observer performs c measurements on each
 203 subject, we extend the two-way random effects to:

$$207 \quad Y_{ijk} = \mu + A_i + B_j + E_{ijk},$$

204 where Y_{ijk} is the k^{th} measurement performed by the j^{th} observer on the i^{th} subject for $i = 1, \dots, a, j = 1, \dots, b,$
 205 and $k = 1, \dots, c$. Note that, conditional on observer and subject, the c repeated measurements are assumed
 206 to be independent and identically distributed.

208 Mimicking the arguments for the single measurement case, but now considering the differences $D_{ijk} =$
 209 $Y_{ijk} - \bar{Y}_{i..}$, we propose the following 95% LOAM:

$$212 \quad \pm 1.96 \sqrt{\frac{b-1}{b} \sigma_B^2 + \frac{bc-1}{bc} \sigma_E^2}.$$

210 Again $\sigma_A^2, \sigma_B^2,$ and σ_E^2 are estimated by the ANOVA estimates (see, e.g., Chapter 4 of Searle et al. [3]), which
 211 are given by

$$213 \quad \hat{\sigma}_A^2 = \frac{MSA - MSE}{bc}, \quad \hat{\sigma}_B^2 = \frac{MSB - MSE}{ac}, \quad \hat{\sigma}_E^2 = MSE,$$

214 where now $MSA = SSA/v_A$, $MSB = SSB/v_B$, and $MSE = SSE/v_E$ with $SSA = bc \sum_{i=1}^a (\bar{y}_{i..} - \bar{y}_{...})^2$, $SSB =$
215 $ac \sum_{j=1}^b (\bar{y}_{.j.} - \bar{y}_{...})^2$, $SSE = \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^c (y_{ijk} - \bar{y}_{i..} - \bar{y}_{.j.} - \bar{y}_{...})^2$, and $v_E = abc - a - b + 1$, while
216 $v_A = a - 1$ and $v_B = b - 1$ is unchanged.

217 Note that the overall, subject-specific, and observer-specific averages ($\bar{y}_{...}$, $\bar{y}_{i..}$, and $\bar{y}_{.j.}$) are now also
218 averaging across the multiple measurement index. With these definitions of SSB , SSE , v_B , and v_E and with
219 $N = abc$, the LOAM estimate and CIs still have the form given by Eq. (4)-**Error! Reference source not**
220 **found..** For the sample size calculation summarised in Eq. (6)-(7), we furthermore replace a with ac .

221 Further, CIs for σ_A , σ_B , and σ_E are obtained by Eq. (8)-(10), except that a is replaced with ac , b is replaced
222 by bc , and the definition of $\hat{\sigma}_A^2$, $\hat{\sigma}_B^2$, $\hat{\sigma}_E^2$, v_A , v_B , and v_E has changed to the above.

223 Note that all formulas for the multiple measurement case reduce to those for the single measurement
224 case, when $c = 1$.

225 As for the single measurement setup, the observations may be visualised using an agreement plot,
226 where the observed differences $d_{ijk} = y_{ijk} - \bar{y}_{i..}$ are plotted against the subject-specific averages $\bar{y}_{i..}$.

227 2.3 Data and software

228 The statistical programming language R, version 3.6.1 [6], was used to analyse the data in the paper. An R-
229 package, R-scripts, and the aortic data for the LOAM calculations in the present paper can be obtained from
230 the GitHub repository: <https://github.com/HaemAalborg/loamr>.

231 3 Results

232 **Example 1.** In a study $b = 5$ thoracic radiologists measured the diameter (in centimetres) of $a = 40$ lung
233 tumours from computed tomography scans [6]. This study was also used as an example in Jones et al. [2].
234 Table 1 shows the empirical mean and standard deviation of the measurements across subject, conditional
235 on radiologist, and Figure 1 displays the agreement plot. Estimates and CIs of the 95% LOAM, ICC, σ_A , σ_B ,
236 and σ_E are listed in Table 2. Neither the agreement plot nor the conditional empirical mean indicate any
237 observer systematically making unusually small or large measurements. Further, there is no indication of
238 heteroscedasticity in relation to change in observer or to the size of the tumour.

239 The estimated 95% LOAM are ± 1.1 cm (95% CI: 1.0 cm to 1.8 cm); the estimate is identical with the 95%
 240 LOAM calculated by Jones et al.'s method when rounding to one decimal place. The inter-observer standard
 241 deviation estimate is 0.3 cm (95% CI: 0.1 cm to 0.5 cm), while the intra-observer standard deviation
 242 estimate is 0.6 cm (95% CI: 0.5 cm to 0.6 cm). Although on a scale comparable to the intra-observer
 243 variation, the inter-observer variation is smaller, supporting the practice where lung nodule measurements
 244 are performed by different radiologists. We may also note that the inter-subject variation (unsurprisingly)
 245 is larger than both the inter- and intra-observer variation.

246

247

<i>Radiologist</i>	<i>Mean (cm)</i>	<i>SD (cm)</i>
1	3.9	1.6
2	3.7	1.5
3	4.4	1.6
4	4.4	1.6
5	4.1	1.6

248

249

250

251

252 Table 1. Empirical mean and standard deviation (SD) of the tumour measurements, calculated across
 253 subjects, conditional on radiologist.

254

255 <Figure 1 here>

256 Figure 1. Agreement plot for tumour size measurements in centimetres with the proposed 95% LOAM
 257 (dashed line) and associated 95% CI (shading).

258

<i>LOAM (CI)</i> <i>in cm</i>	<i>ICC (CI)</i>	$\hat{\sigma}_A$ (CI) <i>in cm</i>	$\hat{\sigma}_B$ (CI) <i>in cm</i>	$\hat{\sigma}_E$ (CI) <i>in cm</i>
1.1 (1.0, 1.8)	0.8 (0.7, 0.90)	1.5 (1.1, 1.8)	0.3 (0.1, 0.5)	0.6 (0.5, 0.7)

259 Table 2. Estimates and 95% confidence intervals (CIs) of the upper 95% LOAM, intra-class correlation
 260 (ICC), σ_A , σ_B , and σ_E for the tumour measurements.

261

262 **Example 2.** Borgbjerg et al. consider three methods (OTO, LTL, and ITI) for assessing the maximum antero-
 263 posterior abdominal aortic diameter [9]. A total of $b = 12$ radiologists measured the aortic diameter $c = 2$
 264 times on $a = 50$ still abdominal aortic images to assess which of the three methods were most reliable.

265 Using the methods described in Section 2.2 for multiple measurements, we calculate estimates and CIs
 266 for the 95% LOAM, σ_A , σ_B , and σ_E (see Table 3) and make an agreement plot (see Figure 2). The inter-subject
 267 variation is large compared to both the inter- and intra-observer variation. The inter-observer variation is
 268 of the same order of magnitude as the intra-observer variation and should not be excluded. The LTL method
 269 has the largest estimated LOAM, meaning that measurements made by this method tend to vary more.
 270 Conversely, the ITI method has the smallest LOAM suggesting that this method has the highest
 271 reproducibility when taking into account both the inter-observer and intra-observer variation. However,
 272 the wide CIs for the LOAM indicate that more observers may be needed to assess this properly. We found
 273 significantly less intra-observer variation for the LTL and ITI compared to the OTO method. This finding is
 274 in line with the conclusion by Borgbjerg et al. which suggests that it is advantageous to employ either the
 275 ITI or LTL method when repeated measurements are performed by the same observer [9].

276

277

<i>Method</i>	<i>LOAM (CI)</i> <i>in mm</i>	$\hat{\sigma}_A$ (CI) <i>in mm</i>	$\hat{\sigma}_B$ (CI) <i>in mm</i>	$\hat{\sigma}_E$ (CI) <i>in mm</i>
OTO	3.2 (2.8, 4.3)	7.2 (5.7, 8.6)	1.1 (0.7, 1.6)	1.2 (1.2, 1.3)
LTL	3.4 (2.8, 5.1)	6.9 (5.5, 8.3)	1.5 (0.8, 2.1)	1.0 (1.0, 1.1)
ITI	2.9 (2.4, 4.3)	6.8 (5.4, 8.1)	1.2 (0.7, 1.8)	0.9 (0.9, 0.9)

278 Table 3. Estimates and 95% confidence intervals (CIs) for the upper 95% LOAM, σ_A , σ_B , and σ_E for the aortic
 279 diameter measurements.

280

281 <Figure 2 here>

282 Figure 2. Agreement plots for each of the three methods (OTO, LTL, and ITI) used to measure the aortic
 283 diameter along with the estimate (dashed line) and the 95% CI for the 95% LOAM (shading).

284

285 4 Discussion

286 In this study, we have defined the LOAM under the assumption of a two-way random effects model, with
287 additive observer and subject effects. This allowed us to formulate a simple statistical inference procedure
288 which can be easily implemented. The theory could be altered to cover various situations where the
289 assumptions of the paper are not fulfilled.

290 First, we include observers as a random effect, meaning that we consider the observers in a study to be
291 a random sample from a larger population of observers that we want to make inference about. It is, however,
292 not unlikely to have a study where the considered observers constitute the whole population of interest, in
293 which case it may be more appropriate to include observers as a fixed effect. The LOAM presented in this
294 paper is based on the variance of the difference between an individual measurement and the subject-
295 specific mean. Under a model with observers as fixed effect, such a LOAM will no longer measure variation
296 due to change of observer. Depending on the purpose of the agreement study, the estimated observer effects
297 could then be included in a reformulation of the LOAM or considered separately. However, we believe that
298 many studies are performed to investigate agreement not only between the specific observers but rather
299 within a larger population of observers, encouraging the choice of model in this paper.

300 Second, one could imagine a situation where it is relevant to include an interaction term between
301 subjects and observers, that is, modelling that observers may react differently upon the subjects. For single
302 measurements this interaction effect is confounded with the residual error, but for multiple measurements
303 this effect could in principle be modelled and the LOAM adjusted accordingly.

304 Third, the methods and formulae of this paper rely on the assumption of a balanced data setup, where
305 all observers have evaluated all the subjects the same number of times. However, in practice it is not unlikely
306 to encounter an unbalanced data set as measurements may get lost or not all observers were able to perform
307 all measurements. An unbalanced setup is definitely more complicated to handle but some advances can be
308 made. A new expression for the LOAM may be found under a two-way random model allowing unbalanced
309 data, while existing methods for finding estimates of the variance components can be used to estimate the
310 adjusted LOAM (see, e.g., [3], [10]). However, it is in general not possible to obtain closed form expressions
311 for the confidence intervals for the LOAM and variance components.

312 Fourth, as indicated in Section 2.1.5 it might happen that the estimate $\hat{\sigma}_B^2$ is negative due to negative
313 correlation between observations made by the same observer on different subjects which will indicate a
314 misspecification of the two-way random effects model formulated in Eq. (1). It is possible to generalise the
315 theory by considering marginal modelling [11]. It was further indicated in Section 2.1.5 that negativity can
316 also arise by sampling variation of the unbiased ANOVA estimates, we have used in this paper. Various
317 approaches have been suggested to remedy this problem as well [12].

318 Pursuing these generalisations will, however, make modelling and implementation much more
319 involved, and thereby violate our goal to formulate an easily implementable framework.

320 **5 Conclusions**

321 Our results show it is possible to formulate measures for the agreement with the mean between multiple
322 observers, equip them with confidence intervals, and extend them to multiple observations per observer,
323 thereby providing a natural extension of Bland-Altman's graphical method. We believe, we have provided
324 an easily accessible and useful statistical toolbox for researchers involved in assessing agreement between
325 methods or individuals performing clinical measurements.

326 **List of abbreviations**

327	LOAM	Limits of agreement with the mean
328	ICC	Inter-class correlation
329	ANOVA	Analysis of variance
330	CI	Confidence interval

331 **Declarations**

332 **Ethics approval and consent to participate**

333 Regarding the ethics approval and consent to participate, we refer to the statements in the original papers
334 by Erasmus et al. [13] for the tumour sizes data and Borgbjerg et al. [9] for the abdominal aortic diameter
335 measurement data.

336 **Consent for publication**

337 Not applicable

338 **Availability of data and materials**

339 The dataset on abdominal aortic diameter measurements supporting the conclusions of this article is
340 available in the *loamr* repository: <https://github.com/HaemAalborg/loamr>. The dataset on tumour sizes is
341 not publicly available but is available from the corresponding author of the original paper on request [13].

342 **Competing interests**

343 Not applicable

344 **Funding**

345 Not applicable

346 **Authors' contributions**

347 MB and JB designed the study. MB and HSC did the statistical modelling and analysed the data. HSC wrote
348 the first version of the manuscript. LB produced figures and organised data and scripts into an R package.
349 All authors read and approved the final manuscript.

350 **Acknowledgements**

351 Not applicable

352 **Additional material**

353 The following additional pdf files are provided:

354 Additional file 1: Derivation of the confidence intervals for the LOAM

355 Additional file 2: Coverage probabilities from a small simulation study

356 Additional file 3: Derivation of confidence intervals for the variance parameters

357 Additional file 4: Formulae after removing the observer effect

358 **References**

359 [1] J. M. Bland and D. G. Altman, "Statistical methods for assessing agreement between two methods of
360 clinical measurement," *Lancet*, vol. 327, no. 8476, pp. 307–310, 1986, doi: 10.1016/S0140-
361 6736(86)90837-8.

362 [2] M. Jones, A. Dobson, and S. O'brian, "A graphical method for assessing agreement with the mean

- 363 between multiple observers using continuous measures," *Int. J. Epidemiol.*, vol. 40, no. 5, pp. 1308–
364 1313, 2011, doi: 10.1093/ije/dyr109.
- 365 [3] S. R. Searle, G. Casella, and C. E. McCulloch, *Variance Components*. Hoboken: John Wiley & Sons, Inc.,
366 1992.
- 367 [4] F. A. Graybill and C.-M. Wang, "Confidence intervals on nonnegative linear combinations of
368 variances," *J. Am. Stat. Assoc.*, vol. 75, no. 372, pp. 869–873, Dec. 1980, doi:
369 10.1080/01621459.1980.10477565.
- 370 [5] J. C. Pinheiro and D. M. Bates, *Mixed-Effects Models in S and S-PLUS*. Springer, New York, 2000.
- 371 [6] C. E. McCulloch, S. R. Searle, and J. M. Neuhaus, *Generalized, Linear, and Mixed Models*, 2. edition.
372 Hoboken, N. J: John Wiley and Sons, 2008.
- 373 [7] K. O. McGraw and S. P. Wong, "Forming inferences about some intraclass correlation coefficients,"
374 *Psychol. Methods*, vol. 1, no. 1, pp. 30–46, 1996, doi: 10.1037/1082-989X.1.1.30.
- 375 [8] R Core Team, "R: A Language and Environment for Statistical Computing." Vienna, Austria, 2019.
- 376 [9] J. Borgbjerg, M. Bøgsted, J. S. Lindholt, C. Behr-Rasmussen, A. Hørlyck, and J. B. Frøkjær, "Superior
377 reproducibility of the leading to leading edge and inner to inner edge methods in the ultrasound
378 assessment of maximum abdominal aortic diameter," *Eur. J. Vasc. Endovasc. Surg.*, vol. 55, no. 2, pp.
379 206–213, 2018, doi: 10.1016/j.ejvs.2017.11.019.
- 380 [10] R. K. Burdick, C. M. Borrer, and D. C. Montgomery, *Design and Analysis of Gauge R&R Studies: Making*
381 *Decisions with Confidence Intervals in Random and Mixed ANOVA Models*,. SIAM, Philadelphia. ASA,
382 Alexandria, VA: ASA-SIAM Series on Statistics and Applied Probability, 2005.
- 383 [11] G. Mohlenberghs and G. Verbeke, "A note on a hierarchical interpretation for negative variance
384 components," *Stat. Modelling*, vol. 11, no. 5, pp. 389–408, doi: 10.1177/1471082X1001100501.
- 385 [12] André I. Khuri, "Designs for Variance Components Estimation: Past and Present," *Int. Stat. Rev.*, vol.
386 68, no. 3, pp. 311–322, doi: 10.1111/j.1751-5823.2000.tb00333.x.
- 387 [13] J. J. Erasmus *et al.*, "Interobserver and intraobserver variability in measurement of non-small-cell
388 carcinoma lung lesions: Implications for assessment of tumor response," *J. Clin. Oncol.*, vol. 21, no.
389 13, pp. 2574–2582, 2003, doi: 10.1200/JCO.2003.01.144.

390

