Dear Author,

Here are the proofs of your article.

- You can submit your corrections **online**, via **e-mail** or by **fax**.

- For **online** submission please insert your corrections in the online correction form. Always indicate the line number to which the correction refers.

- You can also insert your corrections in the proof PDF and **email** the annotated PDF.

- For fax submission, please ensure that your corrections are clearly legible. Use a fine black pen and write the correction in the margin, not too close to the edge of the page.

- Remember to note the **journal title**, **article number**, and **your name** when sending your response via e-mail or fax.

- **Check** the metadata sheet to make sure that the header information, especially author names and the corresponding affiliations are correctly shown.

- **Check** the questions that may have arisen during copy editing and insert your answers/corrections.

- **Check** that the text is complete and that all figures, tables and their legends are included. Also check the accuracy of special characters, equations, and electronic supplementary material if applicable. If necessary refer to the *Edited manuscript*.

- The publication of inaccurate data such as dosages and units can have serious consequences. Please take particular care that all such details are correct.

- Please **do not** make changes that involve only matters of style. We have generally introduced forms that follow the journal's style.
  Substantial changes in content, e.g., new results, corrected values, title and authorship are not allowed without the approval of the responsible editor. In such a case, please contact the Editorial Office and return his/her consent together with the proof.

- If we do not receive your corrections **within 48 hours**, we will send you a reminder.

- Your article will be published **Online First** approximately one week after receipt of your corrected proofs. This is the **official first publication** citable with the DOI. **Further changes are, therefore, not possible.**

- The **printed version** will follow in a forthcoming issue.

**Please note**

After online publication, subscribers (personal/institutional) to this journal will have access to the complete article via the DOI using the URL: http://dx.doi.org/[DOI].
If you would like to know when your article has been published online, take advantage of our free alert service. For registration and further information go to: http://www.link.springer.com.

Due to the electronic nature of the procedure, the manuscript and the original figures will only be returned to you on special request. When you return your corrections, please inform us if you would like to have these documents returned.

# Metadata of the article that will be visualized in OnlineFirst

| | | |
|---|---|---|
| ArticleTitle | A measure of interrater absolute agreement for ordinal categorical data | |
| Article Sub-Title | | |
| Article CopyRight | Springer-Verlag GmbH Germany, part of Springer Nature (This will be the copyright line in the final PDF) | |
| Journal Name | Statistical Methods & Applications | |

| Corresponding Author | | |
|---|---|---|
| | Family Name | **Bove** |
| | Particle | |
| | Given Name | **Giuseppe** |
| | Suffix | |
| | Division | Dipartimento di Scienze della Formazione |
| | Organization | Università "Roma Tre" |
| | Address | via del Castro Pretorio 20, 00185, Roma, Italy |
| | Phone | |
| | Fax | |
| | Email | giuseppe.bove@uniroma3.it |
| | URL | |
| | ORCID | http://orcid.org/0000-0002-2736-5697 |

| Author | | |
|---|---|---|
| | Family Name | **Conti** |
| | Particle | |
| | Given Name | **Pier Luigi** |
| | Suffix | |
| | Division | Dipartimento di Scienze Statistiche |
| | Organization | Università 'La Sapienza" |
| | Address | P.le A. Moro, 5, 00185, Rome, Italy |
| | Phone | |
| | Fax | |
| | Email | pierluigi.conti@uniroma1.it |
| | URL | |
| | ORCID | |

| Author | | |
|---|---|---|
| | Family Name | **Marella** |
| | Particle | |
| | Given Name | **Daniela** |
| | Suffix | |
| | Division | Dipartimento di Scienze della Formazione |
| | Organization | Università "Roma Tre" |
| | Address | via del Castro Pretorio 20, 00185, Roma, Italy |
| | Phone | |
| | Fax | |
| | Email | daniela.marella@uniroma3.it |
| | URL | |

ORCID

| Abstract | A measure of interrater absolute agreement for ordinal scales is proposed capitalizing on the dispersion index for ordinal variables proposed by Giuseppe Leti. The procedure allows to overcome the limits affecting traditional measures of interrater agreement in different fields of application. An unbiased estimator of the proposed measure is introduced and its sampling properties are investigated. In order to construct confidence intervals for interrater absolute agreement both asymptotic results and bootstrapping methods are used and their performance is evaluated. Simulated data are employed to demonstrate the accuracy and practical utility of the new procedure for assessing agreement. Finally, an application to a real case is provided. |
|---|---|

**ORIGINAL PAPER**

# A measure of interrater absolute agreement for ordinal categorical data

**Giuseppe Bove[2]** · **Pier Luigi Conti[1]** · **Daniela Marella[2]**

## Abstract

A measure of interrater absolute agreement for ordinal scales is proposed capitalizing on the dispersion index for ordinal variables proposed by Giuseppe Leti. The procedure allows to overcome the limits affecting traditional measures of interrater agreement in different fields of application. An unbiased estimator of the proposed measure is introduced and its sampling properties are investigated. In order to construct confidence intervals for interrater absolute agreement both asymptotic results and bootstrapping methods are used and their performance is evaluated. Simulated data are employed to demonstrate the accuracy and practical utility of the new procedure for assessing agreement. Finally, an application to a real case is provided.

**Keywords** Ordinal data · Interrater agreement · Resampling

## 1 Introduction

Ordinal rating scales are frequently developed in study designs where several raters (or judges) evaluate a group of targets. For instance, in language studies new rating scales before their routine application are tested out by a group of raters, who assess the language proficiency of a corpus of argumentative (written or oral) texts

✉ Giuseppe Bove
giuseppe.bove@uniroma3.it

Pier Luigi Conti
pierluigi.conti@uniroma1.it

Daniela Marella
daniela.marella@uniroma3.it

[1] Dipartimento di Scienze Statistiche, Università 'La Sapienza", P.le A. Moro, 5, 00185 Rome, Italy

[2] Dipartimento di Scienze della Formazione, Università "Roma Tre", via del Castro Pretorio 20, 00185 Roma, Italy

29 produced by a group of writers. Similar situations can be found in organizational,
30 educational, biomedical, social, and behavioural research areas, where raters can be
31 counsellors, teachers, clinicians, evaluators, or consumers and targets can be
32 organization members, students, patients, subjects, or objects. When each rater
33 evaluates each target, the raters provide comparable categorizations of the targets.
34 The more the raters categorizations coincide, the more the rating scale can be used
35 with confidence without worrying about which raters produced those categoriza-
36 tions. Hence, the main interest here consists in analysing the extent that raters assign
37 the same (or very similar) values on the rating scale (interrater absolute agreement),
38 that is to establish to what extent raters evaluations are close to an equality
39 relationship (*e.g.*, in the case of only two raters, if the two sets of ratings are
40 represented by *x* and *y* the relation of interest is *x* = *y*). Measures of interrater
41 absolute agreement, as Cohen's Kappa [and extensions to take into account three or
42 more raters, e.g., von Eye and Mun (2005)] and intraclass correlations (ICC)
43 [(Shrout and Fleiss 1979; McGraw and Wong 1996)] are usually applied when
44 dealing with rating performed by ordinal scales. A first problem of these procedures
45 is that they are not originally defined for ordinal scales, and so they have to be
46 adapted. For instance, the application of indices based on Cohen's Kappa need to
47 assign numerical values to the ordinal level of the scale; intraclass correlation
48 indices are based on ANOVA for repeated measures approach for interval data.
49 Another limitation of the above mentioned measures is that they are affected by the
50 restriction of variance problem [e.g., LeBreton et al. (2003)], that consists in an
51 attenuation of estimates of rating similarity caused by an artefact reduction of the
52 between-subjects variance in ratings. For instance, this happens in language studies
53 when the same task is defined for native (L1) and non-native (L2) writers, and the
54 analysis compare rater agreement in the two groups separately. Even in the presence
55 of a very good absolute agreement, Cohen's Kappa coefficient and intraclass
56 correlations can take low values, especially for L1 group, because the range of
57 ratings provided by the raters are concentrated on one or two very high levels of the
58 scale (a range restriction that determines a between-target variance restriction).

59 In order to overcome the restriction of variance problem, measure for absolute
60 agreement (or consensus) have been proposed, see (LeBreton and Senter 2008) for a
61 review. The main underlying idea is to measure the within-target variance of ratings
62 (i.e., the between-rater variance) separately for each target, and summarize the
63 results in a final average index (usually normalized in the interval [0, 1]). In this
64 approach, the influence of the low level of the between-target variance is removed
65 by separate analysis of the ratings of each target. One of the most popular index in
66 this group, denoted by $r_{WG}$, was proposed by James et al. (1984), (1993). Let *X* be
67 an ordinal categorical variable with *K* categories (*e.g.* a Likert scale), the index $r_{WG}$
68 can be expressed as

$$r_{WG} = 1 - \frac{s_X^2}{\sigma_E^2} \tag{1}$$

70 where $s_X^2$ is the observed between-rater variance of the ratings and $\sigma_E^2$ is the
71 between-rater variance obtained from a theoretical null distribution representing a

72 complete lack of agreement among raters. Roughly speaking, the null distribution
73 conceptually represents no agreement, which means that to calculate $r_{WG}$, one
74 makes a direct comparison between the observed variance in raters' ratings with the
75 variance one would expect if there was no agreement among raters. Higher numbers
76 indicate a greater agreement.

77 For raters in perfect agreement we have $s_X^2 = 0$, with a corresponding value
78 $r_{WG} = 1$. In applications, $r_{WG}$ values greater than 0.7 (possibly 0.8) are considered
79 associated with high level of interrater absolute agreement [see (LeBreton and
80 Senter 2008), p. 836 Table 2]. Often researchers define the no agreement, or the null
81 distribution, in terms of a uniform distribution. When the null distribution is
82 assumed as uniform, the equation for the corresponding variance is

$$\sigma_E^2 = \frac{K^2 - 1}{12} \tag{2}$$

84 where $K$ refers to the total number of levels of the scale $X$.

85 The index $r_{WG}$ and other indices reviewed in LeBreton and Senter (2008) (*e.g.*,
86 standard and average deviation indices) allow to avoid the problem of variance
87 restriction, but as traditional measures of interrater agreement they are defined only
88 for interval data. Besides, the accuracy of $r_{WG}$ depends strongly on the specification
89 of the null distribution. One disadvantage of $r_{WG}$ is the ambiguity in choosing the
90 reference distribution. Although (James et al. 1984) recommended using the
91 uniform distribution, Lindell and Brandt (1997) recommended using maximum
92 dissensus. Burke et al. (1999), however, cautioned against the use of maximum
93 dissensus because its use may lead to the overestimation of interrater agreement.
94 Finally, depending on the choice of the null distribution, negative values could be
95 obtained for $r_{WG}$. For these reasons, in this contribution we propose a new procedure
96 to measure absolute agreement for ordinal rating scales by using the dispersion
97 index proposed by Leti (1983) (pp. 290–297) for ordinal variables. In this way, we
98 take into consideration the ordinal level of the measurement scales. The new
99 measure is not affected by restriction of variance problems and does not depend on
100 the choice of a particular null distribution. In this paper we assume a two-way
101 random sampling design, where the sampling design involves a sample of raters as
102 well as a sample of targets, all of which are rated by each sampled rater.

103 The paper is organized as follows. In Sect. 2 the dispersion index proposed by
104 Leti (1983) (pp. 290–297) for ordinal variables is introduced and its sampling
105 properties are analyzed in Sect. 3. Such results allow to construct confidence
106 interval without resorting to bootstrap method, as generally happened for inference
107 on measure of interrater absolute agreement, see (Cohen et al. 2001) and reference
108 therein. Section 4 contains results of a simulation experiment used to illustrate both
109 the performance of the proposed interrater agreement index and to compare it with
110 the bootstrap method in constructing confidence intervals. Finally, in Sect. 5 an
111 application to real data is performed.

## 2 Leti index as a measure of interrater absolute agreement for ordinal scales

The dispersion of an ordinal categorical variable can be measured by the index proposed in Leti ([1983](#)) (pp. 290–297), which is given by

$$D = 2 \sum_{k=1}^{K-1} F_k(1 - F_k) \qquad (3)$$

where $K$ is the number of categories of the variable $X$ and $F_k$ is the cumulative proportion associated to category $k$, for $k = 1, \ldots, K$. It is interesting to notice that $D$ has properties of within and between dispersion decomposition analogous to the well-known variance decomposition (Grilli and Rampichini [2002](#)). Index ([3](#)) is nonnegative and it is easy to prove that $D = 0$ if and only if all observed categories are equal (absence of dispersion). The maximum value of the index ($D_{max}$) is obtained when all observations are concentrated in the two extreme categories of the variable (maximum dispersion), and it is

$$D_{max} = \frac{K - 1}{2} \qquad (4)$$

as $N$ is even,

$$D_{max} = \frac{K - 1}{2}\left(1 - \frac{1}{N^2}\right) \qquad (5)$$

as $N$ is odd, $N$ being the total number of observations. For $N$ moderately large, the maximum of the index can be assumed equal to $(K - 1)/2$. Hence, it is possible to define a measure of dispersion normalized in the interval [0, 1] given by

$$d = \frac{D}{D_{max}} = \frac{2}{K - 1}D. \qquad (6)$$

The lower the value of $d$ the higher the raters agreement. Note that, when $d = 0$ (maximum agreement between raters) $r_{WG} = 1$. When $d = 1$

$$r_{WG} = 1 - \frac{(K - 1)^2}{4}\frac{1}{\sigma_E^2} \qquad (7)$$

and if the uniform distribution is assumed as null distribution, ([7](#)) becomes

$$r_{WG} = \frac{4 - 2K}{K + 1} \qquad (8)$$

taking value lower than zero when $K > 2$. In accordance with (LeBreton et al. [2005](#)) out-of-bounds values ($r_{WG} < 0$ or $r_{WG} > 1$) are generally setted to zero. Unlike $r_{WG}$, $d$ can never be out of the range [0, 1].

Advantages of our proposal respect to measures of absolute agreement like $r_{WG}$ are:

142  (i)  *d* takes into consideration the ordinal level of the measurement scales;
143  (ii)  *d* allows to avoid the problem of restriction of variance;
144  (iii)  *d* does not depend by the formulation of a null distribution for
145       normalization;
146  (iv)  the sampling proprieties of *d* are known, as showed in Sect. 3.

148  **Remark 1** In order to homogenize the values assumed by *d* and $r_{WG}$, the index
149  $1 - d$ can be considered.

150  Suggestions for interpreting the value of $1 - d$ appropriately are in Table 1,
151  where a comparison between $r_{WG}$ and *d* $(1 - d)$ is reported. More specifically,
152  datasets with different level of raters agreement have been generated and the indices
153  $r_{WG}$, *d* and $1 - d$ have been computed.
154  As reported in LeBreton et al. (2003) values of $r_{WG}$ greater than 0.7 (possibly
155  0.8) are considered associated with high level of interrater absolute agreement. As
156  shown in Table 1 the same consideration holds for the $1 - d$ index.
157  Finally, in this paper a single item on Likert scale with *K* categories has been
158  considered. For *J* items, the index $r_{WG}$ (denoted by $r_{WG(J)}$) can been defined as
159  shown in Cohen et al. (2001). Analogoulsy to $r_{WG(J)}$, extensions to *J* items for *d*
160  index based on the average of *J* values $d_j$, each computed for each single item, can
161  be considered.

162  ## 3 Sampling properties of *d* index

163  A sample of $n_R$ raters and a sample of $n_T$ targets are drawn by simple random
164  sampling without replacement from a finite population of targets and raters,
165  respectively. Let us denote with $X_{ij}$ the score given by the *j*th rater to the *i*th target
166  on a *K*-point scale, for $i = 1,\ldots,n_T$ and $j = 1,\ldots,n_R$. Formally, $X_{ij}$s are
167  independent categorical random variables having *K* categories with
168  $p_k^{(ij)} = P(X_{ij} = k)$, for $i = 1,\ldots,n_T$, $j = 1,\ldots,n_R$ and $k = 1,\ldots,K$. In the sequel
169  we assume that both the targets and the raters are homogeneus (*targets-raters*
170  *homogeneity assumption*), which implies that the probability $p_k^{(ij)} = p_k$ does not
171  depend on rater *j* or target *i*, for $i = 1,\ldots,n_T$, $j = 1,\ldots,n_R$, $k = 1,\ldots,K$. As a
172  consequence of *homogeneity assumptions*, the variables $X_{ij}$ are independent and
173  identically distributed (*i.i.d.*).

**Table 1** Comparison between $r_{WG}$ and *d* $(1 - d)$

| $r_{WG}$ | $d$ | $1 - d$ |
|---|---|---|
| 0.07 | 0.81 | 0.19 |
| 0.34 | 0.61 | 0.39 |
| 0.49 | 0.53 | 0.47 |
| 0.74 | 0.32 | 0.68 |
| 0.83 | 0.14 | 0.86 |

174 **Remark 2** With regard to the raters homogeneity, variability in scores provided by
175 raters may depend on a number of raters characteristics such as their expertise,
176 familiarity with the assessment process, or amount of training raters received prior
177 to the rating task, etc. Cumming et al. (2002) showed that rating was positively
178 influenced by earlier rating experience and by experience as an EFL/ESL or English
179 L1 teacher. Thompson (1991) indicated that training in linguistics and knowledge of
180 other languages may lead to higher degrees of interrater reliability. Roughly
181 speaking, assuming raters homogeneity means to eliminate the effect of such
182 characteristics on raters score.

183 Evaluations of interrater agreement can be applied to a number of different
184 contexts and are frequently encountered in social, medicine, psychology and
185 education. An application in medicine and in education are illustrated in Examples 1
186 and 2, respectively.

187 **Example 1** Gleason grading is a used grading system for prostatic carcinoma. The
188 Gleason Score is the grading system used to determine the aggressiveness of
189 prostate cancer. This grading system can be used to choose appropriate treatment
190 options. The Gleason Score ranges from 1 to 5 and describes how much the cancer
191 from a biopsy looks like healthy tissue (lower score) or abnormal tissue (higher
192 score). In Allsbrook et al. (2001) 46 needle biopsies containing prostatic carcinoma
193 were assigned Gleason scores by 10 urologic pathologists. Clearly the urologic
194 pathologists do not necessarily give the same grading for each patient. However, we
195 would expect that they tend to agree with each other. The hypothesis that $X_{ij}$ are *i.i.d*
196 comes from the assumption of targets and raters homogeneity. With regard to
197 Allsbrook et al. (2001) study: (i) the 10 urologic pathologists are homogeneous
198 since they have the same background knowledge and familiarity with grading
199 system; (ii) the 46 patients are homogeneous because affected by the same kind of
200 prostatic carcinoma.

201 **Example 2** A study of agreement among raters in educational research is in Kuiken
202 and Vedder (2014), where raters' judgements of writing performance in L2 and L1
203 has been analyzed. More specifically, all texts in L2 and L1 were rated by expert
204 raters on both communicative adequacy and linguistic complexity on a six-point
205 Likert scale. All raters were experienced L2-teachers and native speakers of the
206 target language. Furthermore, they are homogeneous with respect to the familiarity
207 with the assessment process and the amount of training raters received prior to the
208 rating task.

209 As previously stressed, the dispersion of an ordinal categorical variable can be
210 measured by the index (3).

211 With regard to *i*th target, let us denote with $\widehat{F}_k^{(i)}$ the empirical cumulative
212 distribution function defined as

$$\widehat{F}_k^{(i)} = \frac{1}{n_R} \sum_{j=1}^{n_R} I_{(X_{ij} \leq k)} \tag{9}$$

where the numerator represents the number of raters giving score less than or equal to $k$ to the $i$th target. It is known that $E(\widehat{F}_k^{(i)}) = F_k^{(i)} = F_k$, where the last equality comes from the *targets homogeneity assumptions*. Furthermore, $V(\widehat{F}_k^{(i)}) = F_k(1 - F_k)$ and $Cov(\widehat{F}_k^{(i)}, \widehat{F}_l^{(i)}) = min(F_k, F_l) - F_k F_l$. In order to estimate (3), for each target $i$ the following estimator can be defined

$$\widehat{D}_i = 2 \sum_{k=1}^{K-1} \widehat{F}_k^{(i)}(1 - \widehat{F}_k^{(i)}). \tag{10}$$

As stressed in Piccarreta (2001), (10) can be alternatively expressed as

$$\begin{aligned}
\widehat{D}_i &= \sum_{k=1}^{K} \sum_{l=1}^{K} |k - l| \widehat{p}_k^{(i)} \widehat{p}_l^{(i)} \\
&= \frac{1}{n_R^2} \sum_{j=1}^{n_R} \sum_{j'=1}^{n_R} |X_{ij} - X_{ij'}|
\end{aligned} \tag{11}$$

where

$$\widehat{p}_k^{(i)} = \frac{1}{n_R} \sum_{j=1}^{n_R} I_{(X_{ij}=k)} \tag{12}$$

is an unbiased estimator of $p_k$.

**Proposition 1** *The random variable (r.v.) $n_R(\widehat{p}_1, \ldots, \widehat{p}_K)'$, with $\widehat{p}_k = \sum_{i=1}^{n_T} \widehat{p}_k^{(i)}/n_T$ for $k = 1, \ldots, K$, follows a multinomial distribution with parameters $n_R$ and $(p_1, \ldots, p_K)$.*

The expression (11) allows to compute easily the expectation and the variance of estimator (10) as shown in Proposition 2, see Lomnicki (1952) for details.

**Proposition 2** *The estimator $\widehat{D}_i$ has expectation*

$$E(\widehat{D}_i) = \left(1 - \frac{1}{n_R}\right)D \tag{13}$$

*and variance given by*

$$Var(\widehat{D}_i) = \left(\frac{1}{n_R^2} - \frac{1}{n_R^3}\right)(4\sigma^2 + 4(n_R - 2)J - 2(2n_R - 3)D^2) = V \tag{14}$$

*where*

$$\sigma^2 = Var(X_{ij}) = \sum_{k=1}^{K} k^2 p_k - \left( \sum_{k=1}^{K} k p_k \right)^2 \tag{15}$$

236

$$J = \sum_{k=1}^{K} \sum_{h=1}^{K} \sum_{l=1}^{K} |k - h||k - l| p_k p_h p_l. \tag{16}$$

237

238 **Proof** Both (13) and (14) come from the results in Lomnicki (1952). With regard to
239 (13), we have

$$
\begin{aligned}
E(\widehat{D}_i) &= E\left( \frac{1}{n_R^2} \sum_{j=1}^{n_R} \sum_{j'=1}^{n_R} |X_{ij} - X_{ij'}| \right) \\
&= \frac{n_R(n_R - 1)}{n_R^2} E\left( \frac{1}{n_R(n_R - 1)} \sum_{j=1}^{n_R} \sum_{j'=1}^{n_R} |X_{ij} - X_{ij'}| \right) \\
&= \frac{n_R(n_R - 1)}{n_R^2} 2 \sum_{k=1}^{K-1} F_k(1 - F_k) \\
&= \left( \frac{n_R - 1}{n_R} \right) D.
\end{aligned}
\tag{17}
$$

241 For the variance (14) we obtain

$$
\begin{aligned}
Var(\widehat{D}_i) &= Var\left( \frac{1}{n_R^2} \sum_{j=1}^{n_R} \sum_{j'=1}^{n_R} |X_{ij} - X_{ij'}| \right) \\
&= \left( \frac{n_R - 1}{n_R} \right)^2 Var\left( \frac{1}{n_R(n_R - 1)} \sum_{j=1}^{n_R} \sum_{j'=1}^{n_R} |X_{ij} - X_{ij'}| \right) \\
&= \left( \frac{n_R - 1}{n_R} \right)^2 \frac{1}{n_R(n_R - 1)} (4\sigma^2 + 4(n_R - 2)J - 2(2n_R - 3)D^2) \\
&= \left( \frac{1}{n_R^2} - \frac{1}{n_R^3} \right) (4\sigma^2 + 4(n_R - 2)J - 2(2n_R - 3)D^2).
\end{aligned}
\tag{18}
$$

243 □

244 **Remark 3** For $n_R$ sufficiently large, we have

$$Var(\widehat{D}_i) \approx \frac{4(J - D^2)}{n_R}. \tag{19}$$

245

246 As an estimator of $d$ index (6) we consider

$$\widehat{d} = \frac{\overline{\widehat{D}}}{D_{max}} = \frac{1}{D_{max}} \left( \frac{1}{n_T} \sum_{i=1}^{n_T} \widehat{D}_i \right). \tag{20}$$

248 where $\overline{\widehat{D}}$ is an estimator of $D$ obtained averaging the $n_T$ estimates $\widehat{D}_1, \ldots, \widehat{D}_{n_T}$.

249 In Proposition 3 both the sampling properties and the asymptotic distribution of $\widehat{d}$
250 are analyzed for large $n_T$ (*e.g.*, $n_T > 30$) and moderate $n_R$ (*e.g.*, $n_R = 7 - 10$).

251 **Proposition 3** *The estimator $\widehat{d}$ has expectation*

$$E(\widehat{d}) = \left( \frac{n_R - 1}{n_R} \right) d \tag{21}$$

253 *and variance*

$$V_d = \left( \frac{1}{D_{max}} \right)^2 \frac{V}{n_T} \tag{22}$$

255 *where $V$ is given in* (14). *Furthermore, since $\widehat{D}_1, \ldots, \widehat{D}_{n_T}$ are i.i.d., for the central*
256 *limit theorem, as $n_T$ goes to infinity the random variable $\widehat{d}$ tends to a standard*
257 *normal distribution with mean and variance given by* (21) *and* (22),*respectively.*

258 **Remark 4** If the homogeneity assumption is violated then the $X_{ij}$ random variables
259 are independent but not identically distributed. The main result in this area is the
260 Liaponouv's Theorem, (see Billingsley 1995). The theorem strengthens the
261 requirement of finite variance requiring that the $X_{ij}$ have finite moments of order
262 $(2 + \delta)$, for some $\delta > 0$. Clearly, the convergence to normal distribution could be
263 slower.

264 In Proposition 4 an unbiased estimator of $d$ is proposed and its asymptotic
265 distribution is evaluated.

266 **Proposition 4** *From* (21), *an unbiased estimator of d can be defined as follows*

$$\widehat{d^*} = \frac{n_R}{n_R - 1} \widehat{d}. \tag{23}$$

268 *As a consequence of Proposition*(3), *the distribution of $\widehat{d^*}$ is approximately normal*
269 *with mean d and variance*

$$V_{d^*} = \left( \frac{n_R}{n_R - 1} \right)^2 \left( \frac{1}{D_{max}} \right)^2 \frac{V}{n_T}. \tag{24}$$

270

271 The proof of Proposition 4 follows from Proposition 3. The above results are
272 useful to construct point and interval estimates of $d$. They are also useful for testing
273 both the statistical significance of the index (that is the null hypotheses $H_0 : d = 0$)
274 and null hypothesis such as $H_0 : d \leq d_0$, where $d_0$ be a real number in [0, 1].
275 Consider the hypothesis problem

$$\begin{cases} H_0: & d \leq d_0 \\ H_1: & d > d_0 \end{cases} \tag{25}$$

277 As a consequence of Proposition 4, a test with an asymptotic significance level $\alpha$
278 consists in accepting $H_0$ whenever

$$\widehat{d}^* \leq d_0 + z_\alpha \sqrt{\widehat{V}_{d^*}} \tag{26}$$

280 where $z_\alpha$ is the $\alpha$-th quantile of the standard normal distribution and $\widehat{V}_{d^*}$ is an
281 estimate of variance (24).

282    The performance $\widehat{d}^*$ has been evaluated in Sect. 4 by a simulation study and it
283 has been compared with the bootstrap method. With regard to the size of $d$, the
284 judgment depends on the application context. Researchers should gain experience
285 using the proposed index to understand which values might be expected to be
286 obtained for $d$ in various situations and how to interpret these values. For instance,
287 one of the main questions in multilevel data analysis is whether it is appropriate to
288 aggregate data and to use the aggregated measures to make inferences about higher
289 level units. A necessary precondition for aggregation is that there is an agreement
290 among the individuals who form the group with regard to the aggregated construct.
291 In this context, the problem is to evaluate if the degree of agreement justifies data
292 aggregation. From this perspective, the hypothesis test (25) assumes a fundamental
293 importance.

294 **Remark 5** If the index $1 - d$ introduced in Remark 1 is considered, as a
295 consequence of Proposition 4, the distribution of $1 - \widehat{d}^*$ is approximately normal
296 with mean $1 - d$ and variance given by (24).

## 4 Simulation study

298 In this section a simulation study has been performed. The aim is: (i) to evaluate the
299 performance of $\widehat{d}^*$; (ii) to compare the normal approximation for the distribution of
300 $\widehat{d}^*$ with the bootstrap method. Such a method is generally used in constructing
301 confidence intervals of interrater agreement measures but its use is recommended
302 when $n_R$ is sufficiently large (e.g., $n_R > 20$), see Cohen et al. (2001). Alternative
303 methods based on bootstrap to construct confidence intervals are compared in the
304 simulation.

305    We focus on confidence intervals for the index $d$ because confidence intervals
306 indicate the range within which the population parameter $d$ (the interrater agreement
307 in the population) is likely to fall, as well as precision of this estimate (i.e., the size
308 of the range).

309    A finite population of size $N_T = 150$ targets and $N_R = 28$ raters was generated
310 from a multinomial model with parameters $N_R = 28$ and probabilities
311 $(p_1, p_2, p_3, p_4, p_5) = (0.1, 0.2, 0.35, 0.25, 0.1)$. Then, the finite population consists
312 in a matrix $P$ of size $N_T \times N_R$. The value of $d$ index (6) is 0.61.

313 From the population, $S = 1000$ samples were drawn according to a simple
314 random sampling without replacement on the basis of the following two-step
315 procedure. First of all, a simple random sample of size $n_R = 7$ from the $N_R = 28$
316 raters has been selected. This is equivalent to select a simple random sampling
317 without replacement of columns in the finite population matrix $P$, the result is a
318 matrix $P_R$ of size $N_T \times n_R$. Secondly, a simple random sampling of size $n_T = 50$
319 from $N_T = 150$ targets has been drawn. This means to draw a simple random
320 sampling of $n_T = 50$ rows from $P_R$.

321 In order to construct confidence intervals for the index $d$, both the asymptotic
322 result in Proposition 4 and bootstrapping procedures are used. The bootstrap
323 methods are described in points (2)–(4) below, where we assume that $B = 1000$
324 bootstrap samples are drawn from each initial sample $s$. Formally, confidence
325 intervals for $d$ of level $1 - \alpha = 0.95$ have been constructed using the following
326 methods:

327 (1) *Normal approximation* For the initial sample $s$ (for $s = 1, \ldots, S$), the
328 confidence interval $[L^s_{Norm}, U^s_{Norm}]$ based on the asymptotic normal approxi-
329 mation is given by

$$L^s_{Norm} = \widehat{d}^* - z_{1-\alpha/2}\sqrt{\widehat{V}_{d^*}}; \quad U^s_{Norm} = \widehat{d}^* - z_{\alpha/2}\sqrt{\widehat{V}_{d^*}} \qquad (27)$$

331 where $\widehat{d}^*$ and $\widehat{V}_{d^*}$ are the estimates of $d$ and $V_{d^*}$, respectively.
332

333 (2) *Percentile method* For the initial sample $s$ (for $s = 1, \ldots, S$), the confidence
334 interval $[L^s_{Perc}, U^s_{Perc}]$ is obtained by taking $\alpha/2$ and $1 - \alpha/2$ quantiles of the $B$
335 bootstrap samples. Formally

$$L^s_{Perc} = Q_{\alpha/2}; \quad U^s_{Perc} = Q_{1-\alpha/2} \qquad (28)$$

337

339 (3) *Bootstrap-t interval* For the initial sample $s$ (for $s = 1, \ldots, S$), the confidence
340 interval $[L^s_T, U^s_T]$ is computed as follows

$$L^s_{T-int} = \widehat{d^*} - t_{1-\alpha/2}\sqrt{\widehat{V}_{d^*}}; \quad U^s_{T-int} = \widehat{d^*} - t_{\alpha/2}\sqrt{\widehat{V}_{d^*}} \qquad (29)$$

342 where $t_\alpha$ is the $\alpha$th percentile of the distribution of $z^*_b$ (for $b = 1, \ldots, B$) with

$$z^*_b = \frac{\widehat{d^*_b} - \widehat{d^*}}{\widehat{se}^*_b}. \qquad (30)$$

345 In (30) $\widehat{d^*_b}$ is the estimate of $d^*$ based on the $b$th bootstrap sample and $\widehat{se}^*_b$ is
347 the standard error based on the data in the $b$th bootstrap sample.

348 (4) *Pivotal method* For the initial sample $s$ (for $s = 1, \ldots, S$), the confidence
349 interval $[L^s_{Pivot}, U^s_{Pivot}]$ is computed as follows

$$L^s_{Pivot} = 2\widehat{d^*} - Q_{1-\alpha/2}; \quad U^s_{Pivot} = 2\widehat{d^*} - Q_{\alpha/2} \qquad (31)$$

351 where $Q_{\alpha/2}$ and $Q_{1-\alpha/2}$ are the $\alpha/2$ and $1 - \alpha/2$ quantiles of the $B$ bootstrap
353 estimates $\widehat{d^*_b}$, for $b = 1, \ldots, B$.

354 As far as the methods described in steps (2)–(4) are concerned, from each of the
355 $S = 1000$ initial samples, the $B = 1000$ bootstrap samples were selected according
356 to the following methods:

357 1 *Nonparametric bootstrap* From each initial sample $s$, the $b$th bootstrap sample is
358   selected as follows: (i) a simple random sample with replacement of $r = 7$ raters
359   has been selected from the original sample of raters; (ii) a simple random
360   sampling with replacement of $n = 50$ writers has been drawn from the original
361   sample of writers. Then, bootstrap is applied to the raters sample as well as the
362   targets sample in order to take into account the variability in $\widehat{d}^*$ due to the two-
363   way random sampling design (where the sampling design involves a sample of
364   raters and a sample of targets). Clearly, when the sampling design involves only
365   the raters the proposed methodology resembles that used in literature, see Cohen
366   et al. (2001) and reference therein.
367 2 *Parametric bootstrap* From each initial sample $s$, the $b$th bootstrap sample is
368   generated according the multinomial model specified in Proposition 1.
369 3 *Pseudo-Nonparametric bootstrap* The nonparametric bootstrap described in
370   point (1), is based on the assumption that the data are *i.i.d.*, see Efron (1979).
371   Since survey data are not necessarily *i.i.d.*, many bootstrap resampling methods
372   have been proposed in the context of survey sampling. These methods are
373   obtained after making some modifications to the classical *i.i.d.* bootstrap in
374   order to adapt it for survey data. For a review of bootstrap methods in the
375   context of survey data, see Mashreghi et al. (2016). The class of pseudo-
376   population bootstrap methods consists in creating a pseudo-population by
377   repeating the units of the initial sample and drawing from such a pseudo-
378   population bootstrap samples with the same design as the initial one. In order to
379   illustrate how a pseudo-population is constructed, let us assume that a simple
380   random sample without replacement has been selected from a finite population
381   of size $N$. A pseudo-population of size $N$ can be created by repeating the
382   selected sample, $N/n$ times. This method, was first introduced by Gross (1980).
383   In practice $N/n$ is rarely an integer, in this case a method to build a pseudo-
384   population of size $N$ was proposed by Booth et al. (1994). In this method, a
385   pseudo-population is first constructed by replicating $k = \lfloor N/n \rfloor$ times each unit
386   of the original sample $s$. Then, the pseudo-population is completed by taking a
387   simple random sample of size $N - nk$ without replacement from $s$. Taking into
388   account the two-way sampling design of both targets and raters, the pseudo-
389   population has been generated according the following two step procedure:

391 Step 1    the ratings of $N_R = 28$ raters have been reconstructed replicating the
392          columns of the original sample $s$, $k_R = N_R/n_R = 28/7 = 4$ times. As a
393          consequence, this first step generates a sample $s_R$ of size $n_T = 50$ and
394          $n_R = N_R = 28$;
395 Step 2    the points of $N_T = 150$ targets have been reconstructed replicating the rows
396          of the sample $s_R$ obtained in Step 1, $k_T = N_T/n_T = 150/50 = 3$ time.

397 The accuracy of confidence intervals has been evaluated by the following indicators.

398 (1) Estimated coverage probability, in per cent, for the interval

$$ECP = \frac{100}{S}\sum_{s=1}^{S} I(L_t^s \le d \le U_t^s). \tag{32}$$

402 (2) Estimated left-tail and right-tail errors (lower and upper error rates) in per
403 cent

$$LE = \frac{100}{S}\sum_{s=1}^{S} I(L_t^s > d), \tag{33}$$

405

$$RE = \frac{100}{S}\sum_{s=1}^{S} I(U_t^s < d). \tag{34}$$

409 (3) Estimated average length (AL) of all 1000 simulated intervals given by

$$AL = \sum_{s=1}^{S} \frac{U_t^s - L_t^s}{S} \tag{35}$$

413 where $I(a) = 1$ if $a$ is true and $I(a) = 0$ elsewhere, and $t = Norm$,
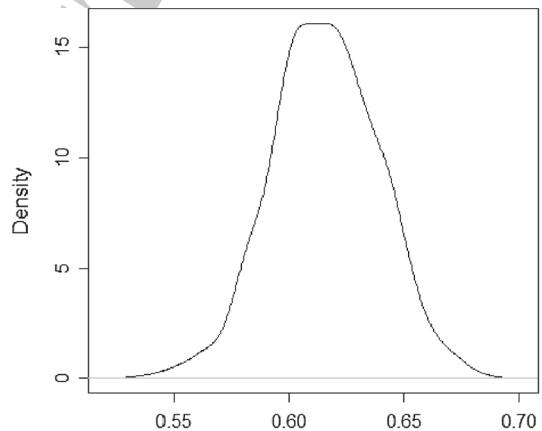414 $T - int, Perc, Pivot$.

### 4.1 Simulation results

416 Tables 2 presents the outcomes achieved in the simulation study. More specifically,
417 the estimated coverage probabilities of 95% confidence intervals (CP), the estimated
418 left-tail (LE) and right-tail (RE) errors (nominal values is 2.5% for both) and the
419 average length (AL) for the index $d$, when $(n_R = 7, n_T = 50)$, are reported. The
420 $d$ value is equal to 0.61.

421 As reported in Table 2, the confidence intervals obtained with the normal
422 approximation perform very well. Coverage probabilities are larger than 95%
423 nominal value (99.4%) with an average length of 0.16. Furthermore, the normal
424 confidence intervals construction is simple, as it does not require resampling from
425 the initial sample. Figure 1 shows the kernel density of the $d$ index estimated from
426 the 1000 original samples. The bandwidth selection rule is as proposed by Sheather
427 and Jones (1991).

428 The percentile method has a good performance with coverage probability larger
429 than 91%. The worst methods are the *Pivot* and *T*-int methods. The lower and upper
430 error rates, giving us an idea of how skewed the distribution of the $d$ estimator is, are
431 not well balanced. With regard to the methods used to generate the bootstrat
432 samples, the *parametric* approach performance is strictly related to the estimation of
433 the multinomial probabilities. As previously stressed, each row in the inital sample $s$
434 provides an estimate of $(p_1, p_2, p_3, p_4, p_5)$ and the mean of such estimates defines the

**Table 2** Performance of different confidence intervals for d when $n_R = 7$, $d = 0.61$

| Method | Indicators | $n_R = 7$ | | |
| --- | --- | --- | --- | --- |
| | | Nonparametric | Parametric | Pseudo-Nonparametric |
| Normal | CP | 99.4 | 99.4 | 99.4 |
| | LE | 0.6 | 0.6 | 0.6 |
| | RE | 0 | 0 | 0 |
| | AL | 0.16 | 0.16 | 0.16 |
| T-int | CP | 26.2 | 72.4 | 28.8 |
| | LE | 73.8 | 26.2 | 71.2 |
| | RE | 0 | 1.4 | 0 |
| | AL | 0.18 | 0.08 | 0.15 |
| Perc | CP | 92.8 | 91.2 | 92.8 |
| | LE | 0 | 8.8 | 0 |
| | RE | 7.2 | 0 | 7.2 |
| | AL | 0.23 | 0.10 | 0.18 |
| Pivot | CP | 27 | 79.2 | 30 |
| | LE | 73 | 19.6 | 70 |
| | RE | 0 | 1.2 | 0 |
| | AL | 0.23 | 0.10 | 0.18 |



**Fig. 1** Kernel density estimate of $d$ index from the 1000 original samples

estimated probabilities $(\widehat{p}_1, \widehat{p}_2, \widehat{p}_3, \widehat{p}_4, \widehat{p}_5)$ of the multinomial distribution used to generate the bootstrap samples as specified in Proposition 1. In Table 3, the minimum, the maximum, the mean and the standard deviation of the distribution of $\widehat{p}_k$ (for $k = 1, 2, 3, 4, 5$) estimated from the original 1000 samples are reported.

**Table 3** Descriptive statistics of $\widehat{p}_k$ distribution, for $k = 1, 2, 3, 4, 5$ and $d = 0.61$

| Parameter | True value | Min | Max | Mean | Sd |
|---|---|---|---|---|---|
| $p_1$ | 0.10 | 0.05 | 0.15 | 0.10 | 0.01 |
| $p_2$ | 0.20 | 0.14 | 0.25 | 0.19 | 0.02 |
| $p_3$ | 0.35 | 0.29 | 0.41 | 0.35 | 0.02 |
| $p_4$ | 0.25 | 0.20 | 0.33 | 0.26 | 0.02 |
| $p_5$ | 0.10 | 0.06 | 0.16 | 0.10 | 0.02 |

439    As Table 2 shows, the *pseudo-nonparametric* approach taking into account the
440 sample selection effects has a slightly better performance than the *nonparametric*
441 approach both in terms of coverage probabilities and average lengths for all methods
442 ($T - int$, *Perc*, *Pivot*).

443    Finally, note that in the *nonparametric* approach the resampling with replace-
444 ment from $n_R = 7$ raters generates a replication of columns of the bootstrap sample
445 introducing a false agreement between raters and as a consequence an underesti-
446 mation of *d*. This fact is showed in Table 4 where the mean of the *d* estimates over
447 both the 1000 original samples *s* and over the bootstrap replications *b* are reported.

448    Such means have been computed both for the original population with $d = 0.61$
449 and for a population with $d = 0.41$, showing as the magnitude of bias depends also
450 on the original agreement degree between raters. That is, the higher the raters
451 agreement (low values of *d*), the smaller the bias in the *d* estimator introduced by
452 the resampling with replacement. Clearly, such a bias is also present in the *pseudo-*
453 *nonparametric* approach but with a smaller magnitude, thank to the construction of
454 the pseudo-population that mitigates such a phenomenon. As Table 4 shows, the
455 *parametric approach* produces null bias estimates.

456    The simulation in Table 2 has been repeated for a population with $d = 0.41$. The
457 results are reported in Table 5.

458    In conclusion, the most competitive method in terms of performance and
459 computational time seem to be the normal. Finally, among the alternative methods
460 based on bootstrap the percentile method in the *parametric* approach seems to
461 perform better.

**Table 4** The mean of $\widehat{d}$ over the initial samples *s* and over the bootstrap replications *b*

| Approach | Mean of $\widehat{d}^*$ *(d=0.61)* | Mean of $\widehat{d}^*$ *(d=0.41)* |
|---|---|---|
| Nonparametric | 0.53 | 0.36 |
| Parametric | 0.61 | 0.41 |
| Pseudo-nonparametric | 0.55 | 0.37 |

**Table 5** Performance of different confidence intervals for d when $n_R = 7$, $d = 0.41$

| Method | Indicators | $n_R = 7$ | | |
| | | Nonparametric | Parametric | Pseudo-nonparametric |
| --- | --- | --- | --- | --- |
| Normal | CP | 98.2 | 98.2 | 98.2 |
| | LE | 1.8 | 1.8 | 1.8 |
| | RE | 0 | 0 | 0 |
| | AL | 0.13 | 0.13 | 0.13 |
| T-int | CP | 60.2 | 83.2 | 61.2 |
| | LE | 39.8 | 14.8 | 38.8 |
| | RE | 0 | 2 | 0 |
| | AL | 0.18 | 0.10 | 0.14 |
| Perc | CP | 93.2 | 93.8 | 93.2 |
| | LE | 0 | 5.8 | 0 |
| | RE | 6.8 | 0.4 | 6.3 |
| | AL | 0.19 | 0.10 | 0.15 |
| Pivot | CP | 64.8 | 84.6 | 65.4 |
| | LE | 35.2 | 12.6 | 34.6 |
| | RE | 0 | 2.8 | 0 |
| | AL | 0.19 | 0.10 | 0.15 |

## 5 An application on real data: the assessment of language proficiency

The aim of this section is to apply the methodology illustrated in the previous sections on an empirical data set, we have analysed ratings obtained in a research conducted at Roma Tre University [see (Nuzzo and Bove 2020), for a detailed description]. The main aim of the study was to investigate the applicability of a six-point Likert scale for functional adequacy (an aspect of language proficiency) developed by Kuiken and Vedder (2017) to texts produced by native and non-native writers, and to different task types (narrative, instruction, and decision-making tasks). The scale comprises four subscales, corresponding to the four dimensions of functional adequacy identified by the authors of the scale: content, task require-ments, comprehensibility, coherence and cohesion [the reader is referred to Kuiken and Vedder (2017) for a detailed presentation of scales and descriptors]. 20 native speakers of Italian (L1) and 20 non-native speakers of Italian (L2) participated in the study as writers. All the texts produced by L1 and L2 writers (120 texts in total for the three tasks) were assessed by 7 native speakers of Italian on the Kuiken and Vedder six-point Likert scale. The raters did not have any specific experience in judging written texts, and can therefore be categorized as being non-expert. For our purposes, we have selected ratings concerning only the narrative task and the subscale comprehensibility. Just to give a general idea of the subscale, definitions of levels 1 and 6 are reported in the following:

483  Level 1:  The text is not at all comprehensible. Ideas and purposes are unclearly
484          stated and the efforts of the reader to understand the text are
485          ineffective.
486  Level 6:  The text is very easily comprehensible and highly readable. The ideas
487          and the purpose are clearly stated.

488  The results of the interrater agreement analysis for the subscale are summarized in
489  Table 6, where the intraclass correlation $ICC(A, 1)$, as defined in McGraw and
490  Wong (1996), and the average values of $r_{WG}$, as defined in LeBreton and Senter
491  (2008), the coefficient of variation $CV$, $\hat{d}$ and $\hat{d}^*$ are shown for L1, L2 and total
492  groups. The intraclass correlation $ICC(A, 1)$ provides a low-moderate level of
493  agreement for the total group (0.67). The results for the average values of $CV$
494  (12.16%) seems in accord with $ICC(A, 1)$, while the average value of $r_{WG} = 0.87$,
495  $\hat{d} = 0.22$ $(1 - \hat{d} = 0.78)$ and $\hat{d}^* = 0.25$ $(1 - \hat{d}^* = 0.75)$ highlight a higher level of
496  agreement. As it was observed in Bove et al. (2018), when the analysis focuses
497  separately on the two subgroups of L1 and L2 students, results regarding the L1
498  group deserve particular attention. Interrater agreement measured by intraclass
499  correlation is very low in the L1 group $(ICC(A, 1) = 0.14)$. Analysing the
500  dispersion of the ratings given to this subgroup, it comes out that most of the
501  raters used almost exclusively levels 5 and 6 of the scale. Such a range restriction
502  caused the very low value of the intraclass correlation, despite the substantial
503  agreement among the raters that scored all the L1 texts in the same high levels. This
504  problem does not regard the results for the other indices of Table 6: $r_{WG} = 0.90$;
505  $CV = 8.12\%$; $\hat{d} = 0.17$ $(1 - \hat{d} = 0.83)$; $\hat{d}^* = 0.19$ $(1 - \hat{d}^* = 0.81)$. that show a
506  very good level of absolute agreement. Finally, the standard deviation of $\hat{d}^*$
507  computed on the basis of formula (24) is equal to 0.05. As a consequence, the
508  $(1 - \alpha) = 0.95$ confidence interval using the normal approximation for the total
509  group is [0.15, 0.35] and the error is at most 0.10.

## 6 Conclusions

511  In this paper a measure of interrater absolute agreement for ordinal scales is
512  proposed. Such a measure is not affected by restriction of variance problems and
513  does not depend on the choice of a particular null distribution. An unbiased
514  estimator of the proposed measure is introduced and its sampling properties are
515  investigated. In the simulation study confidence intervals for the proposed interrater
516  agreement index are constructed using the normal approximation, the parametric

**Table 6** $ICC(A, 1)$ and average of $r_{WG}$, $CV$, $\hat{d}$ and $\hat{d}^*$ for the comprehensibility subscale in the L1, L2 and the total groups

| Group | N | ICC(A, 1) | $r_{WG}$ | CV% | $\hat{d}$ | $\hat{d}^*$ |
|-------|---|-----------|----------|-----|-----------|-------------|
| L1    | 20 | 0.14 | 0.90 | 8.12 | 0.17 | 0.19 |
| L2    | 20 | 0.63 | 0.84 | 16.20 | 0.28 | 0.32 |
| Total | 40 | 0.67 | 0.87 | 12.16 | 0.22 | 0.25 |

🖄 Springer

and nonparametric bootstrap. Furthermore, a pseudo-nonparametric bootstrap taking into account the sampling design is also implemented. As previously stressed, the resampling involves both raters and targets sample. Confidence intervals obtained with the normal approximation seem to perform very well both in terms of coverage probability and computational cost.

# References

Allsbrook WC, Mangold KA, Johnson MH, Lane RB, Lane CG, Amin MB (2001) Interobserver reproducibility of Gleason grading of prostatic carcinoma: urologic pathologists. Hum Pathol 32(1):74–80

Billingsley P (1995) Probability and measure, 3rd edn. Wiley, New York

Booth JG, Butler RW, Hall P (1994) Bootstrap methods for finite populations. J Am Stat Assoc 89(428):1282–1289

Bove G, Nuzzo E, Serafini A (2018) Measurement of interrater agreement for the assessment of language proficiency. In: Capecchi S, Di Iorio F, Simone R. ASMOD 2018: proceedings of the advanced statistical modelling for ordinal data conference. Università Federico II di Napoli. FedOAPress, Napoli pp 61–68

Burke MJ, Finkelstein LM, Dusig MS (1999) On average deviation indices for estimating interrater agreement. Organ Res Methods 2:49–68

Cohen A, Doveh E, Eick U (2001) Statistical properties of the $r_{wg}$ index of agreement. Psychol Methods 6(3):297–310

Cumming A, Kantor R, Powers DE (2002) Decision making while rating ESL/EFL writing tasks: a descriptive framework. Mod Lang J 86:67–96

Efron B (1979) Bootstrap methods: another look at the jackknife. Ann Stat 7(1):1–26

Grilli L, Rampichini C (2002) Scomposizione della dispersione per variabili statistiche ordinali [Dispersion decomposition for ordinal variables]. Statistica 62:111–116

Gross S (1980). Median estimation in sample surveys. In: Proceedings of the section on survey research methods. American Statistical Association, pp. 181–184

James LJ, Demaree RG, Wolf G (1984) Estimating within-group interrater reliability with and without response bias. J Appl Psychol 69:85–98

James LJ, Demaree RG, Wolf G (1993) rwg: an assessment of within-group interrater agreement. J Appl Psychol 78:306–309

James LR, Demaree RG, Wolf G (1984) Estimating within-group interrater reliability with and without response bias. J Appl Psychol 69:85–98

Kuiken F, Vedder I (2014) Rating written performance: What do raters do and why? Lang Test 31(3):329–348

Kuiken F, Vedder I (2017) Functional adequacy in L2 writing: towards a new rating scale. Lang Test 34:321–336

LeBreton JM, Burgess JRD, Kaiser RB, Atchley EK, James LR (2003) The restriction of variance hypothesis and interrater reliability and agreement: are ratings from multiple sources really dissimilar? Organ Res Methods 6:80–128

LeBreton JM, James LR, Lindell MK (2005) Recent issues regarding $rwg$, $r^*wg$, $rwg(j)$, and $r^*wg(j)$. Organ Res Methods 8(1):128–138

LeBreton JM, Senter JL (2008) Answers to 20 questions about interrater reliability and interrater agreement. Organ Res Methods 11(4):815–852

Leti G (1983) Statistica descrittiva. Il Mulino, Bologna

Lindell MK, Brandt CJ (1997) Measuring interrater agreement for ratings of a single target. Appl Psychol Meas 21:271–278

Lomnicki ZA (1952) The standard error of Gini's mean difference. Ann Math Stat 23(14):635–637

Author Proof

568  Mashreghi Z, Haziza D, Léger C (2016) A survey of bootstrap methods in finite population sampling.
569      Stati Surv 10:1–52
570  McGraw KO, Wong SP (1996) Forming inferences about some intraclass correlation coefficients. Psychol
571      Methods 1:30–46
572  Nuzzo E, Bove G (2020) Assessing functional adequacy across tasks: a comparison of learners and native
573      speakers' written texts. EuroAm J Appl Linguist Lang, 2. In print
574  Piccarreta R (2001) A new measure of nominal-ordinal association. J Appl Stat 28(1):107–120
575  Sheather SJ, Jones MC (1991) A reliable data-based bandwidth selection method for Kernel density
576      estimation. J R Stat Soc Ser B 53:683–690
577  Shrout PE, Fleiss JL (1979) Intraclass correlations: uses in assessing reliability. Psychol Bull 86:420–428
578  Thompson I (1991) Foreign accents revisited: factors relating to transfer of accent from the first language
579      to a second language. Lang Speech 24(3):265–272
580  von Eye A, Mun EY (2005) Analyzing rater agreement. Manifest variable methods. Lawrence Erlbaum
581      Associates, Lawrence Erlbaum Associates, Mahwah, New Jersey
582

UNCORRECTED PROOF

🖄 Springer

Author Proof

# Author Query Form

## Please ensure you fill out your response to the queries raised below and return this form along with your corrections

Dear Author

During the process of typesetting your article, the following queries have arisen. Please check your typeset proof carefully against the queries listed below and mark the necessary changes either directly on the proof/online grid or in the 'Author's response' area provided below

| Query | Details Required | Author's Response |
|---|---|---|
| AQ1 | Please check and confirm the author given name and family name is correct. Author [2] Given name [Pier Luigi] Family name [Conti]. Also, kindly confirm the details in the metadata are correct. | |