

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20

Interrater Reliability Estimators Tested against True Interrater Reliabilities

Xinshu Zhao^{*a} (xszhao@um.edu.mo)

Guangchao Charles Feng^a (charlesfeng@um.edu.mo)

Song Harris Ao^a (harrisao@um.edu.mo)

Piper Liping Liu^a (lipingliu@um.edu.mo)

^a: Department of Communication, Faculty of Social Sciences, University of Macau, Taipa,
Macao

Corresponding Author: Xinshu Zhao

21

Abstract

22 **Background:** Interrater reliability, aka intercoder reliability, is defined as true agreement
23 between raters, aka coders, without chance agreement. It is used across many disciplines
24 including medical and health research to measure the quality of ratings, coding, diagnoses, or
25 other observations and judgements. While numerous indices of interrater reliability are
26 available, experts disagree on which ones are legitimate or more appropriate.

27 Almost all agree that percent agreement (a_o), the oldest and the simplest index, is also
28 the most flawed because it fails to estimate and remove chance agreement, which is produced
29 by raters' random rating. The experts, however, disagree on which of the chance-adjusted
30 indices are legitimate or better. The experts also disagree on which of the three factors, rating
31 category, distribution skew, or task difficulty, a good index should rely on to estimate chance
32 agreement, or which of the factors the known indices in fact rely on.

33 The most popular chance-adjusted indices, according to a functionalist view of
34 mathematical statistics, assume that all raters conduct intentional and maximum random
35 rating while typical raters conduct involuntary and reluctant random rating. The mismatch
36 between the assumed and the actual rater behaviors causes the indices to rely on mistaken
37 factors to estimate chance agreement, leading to the numerous paradoxes, abnormalities, and
38 other misbehaviors of the indices identified by prior studies.

39 **Methods:** We conducted a 4×8×3 between-subject controlled experiment with 4 subjects per
40 cell. Each subject was a rating session with 100 pairs of rating by two raters, totaling 384
41 rating sessions as the experimental subjects. The experiment tested seven best-known indices
42 of interrater reliability against the observed reliability and chance agreement. Impacts of the
43 three factors, i.e., rating category, distribution skew, and task difficulty, on the indices were
44 tested.

45 **Results:** The most criticized index, percent agreement (a_o), showed as the most accurate
46 predictor of reliability, reporting directional $r^2=.84$. It was also the third best approximator,
47 overestimating observed reliability by 13 percentage points. The three most acclaimed and
48 most popular indices, Scott's π , Cohen's κ and Krippendorff's α , underperformed all other
49 indices, reporting directional $r^2=.312$ and underestimated reliability by 31.4~31.8 points. The
50 newest index, Gwet's AC_1 , emerged as the second-best **predictor** and the most accurate
51 approximator. Bennett et al's S ranked behind AC_1 , and Perreault and Leigh's I_r ranked the
52 fourth both for prediction and approximation. The reliance on category and skew and failure
53 to rely on difficulty explain why the six chance-adjusted indices often underperformed a_o ,
54 which they were created to outperform. The evidence corroborated the notion that the chance-
55 adjusted indices assume intentional and maximum random rating by the raters while the
56 raters instead exhibited involuntary and unwilling random rating.

57 **Conclusion:** The authors call for more empirical studies and especially more controlled
58 experiments to falsify or qualify this study. If the main findings are replicated and the
59 underlying theories supported, new thinking and new indices may be needed. Index designers
60 may need to refrain from assuming intentional and maximum random rating, and instead
61 assume involuntary and reluctant random rating. Accordingly, the new indices may need to
62 rely on task difficulty, rather than distribution skew or rating category, to estimate chance
63 agreement.

64 **Key words:** intercoder reliability, interrater reliability, reconstructed experiment, Cohen's
65 *kappa*, Krippendorff's *alpha*,.

66

67 **Interrater Reliability Estimators Tested against True Interrater Reliabilities**

68

69 **Background**

70 *Intercoder or interrater reliability* is used to measure measurement quality in many
71 disciplines, including health and medical research (1–10). A search of databases including
72 Google Scholar, Scopus, and Web of Science found dozens of terms in academic literature,
73 such as *diagnostician* for inter-diagnostician reliability and *patient* for inter-patient reliability,
74 showing the concept’s broad reach --

75 annotator, arbitrator, assessor, auditor, diagnostician, doctor, editor, evaluator,
76 examiner, grader, interpreter, interviewer, judge, monitor, observer, operator, patient,
77 pharmacist, physician, reader, referee, reporter, researcher, respondent, scorer,
78 screener, student, supervisor, surgeon, teacher, tester, therapist, transcriber, translator,
79 user, voter.

80 Likely the earliest index is *percent agreement*, denoted a_o (9,11). Almost all reliability
81 experts agree that a_o inflates reliability because it fails to remove *chance agreement* (a_c) (2–
82 5,12–14). Scores of indices have been proposed to estimate and remove a_c . Bennett and
83 colleagues’ S and Perreault and Leigh’s I_r estimate a_c as functions of *category* (C) (7,15).
84 Scott’s π , Cohen’s κ and Krippendorff’s α estimate a_c as functions of distribution *skew* (s_k)
85 (2,16–19). Gwet’s AC_1 makes a_c a function of both category and skew. Although many other
86 indices are available and new indices continue to emerge, only these seven are in regular use
87 and continue to be recommended or advocated, according to comprehensive reviews (14,20–
88 26).

89 Using derivation or simulation, statisticians discuss and debate three questions: 1)
90 Which indices are valid or more accurate when estimating reliability or chance agreement? 2)
91 What factors affect the indices? 3) What factors should affect the indices? Answers to
92 Questions 2 and 3 explain the answers to Question 1 (14,27). Underlying the debates are five
93 viewpoints, the first of which is widely shared by almost all experts, while the others are
94 contested, often heatedly. The five viewpoints lead to five groups of conjectures, which we
95 list below and leave the details to Appendix, Section I.2.

- 96 1. Percent agreement (a_o) ignores chance agreement (a_c), therefore is inflated.
- 97 2. Rating category (C) inflates S , I_r , and AC_I by deflating the indices' a_c
98 estimates.
- 99 3. Distribution skew (s_k) deflates π , κ & α by inflating the indices' a_c estimates.
- 100 4. Major indices overlook task difficulty, a major factor affecting a_c ;
101 consequently, they misestimate reliability.
- 102 5. Chance-adjusted indices, S , π , κ , α , I_r , and AC_I included, assume intentional
103 and maximum chance rating by all raters; it is under this assumption that the
104 chance-adjusted indices share the same chance correcting formula, Equation 1,
105 where a_o is observed %-agreement, a_c is estimated chance agreement, and r_i is
106 estimated true agreement, i.e., reliability index.

$$r_i = \frac{a_o - a_c}{1 - a_c} \quad (1)$$

107 The intentional-random assumption, aka maximum-random assumption, is said to be a
108 root cause of many known paradoxes, abnormalities, and other misbehaviors of the indices,
109 because raters are believed to be honest and truthful. Random ratings, if any, should be
110 involuntary rather than intentional, task-dependent rather than invariably maximized (14,21–
111 24,26,28–30).

112 Chance agreement is a product of rater behavior, and the debates are ultimately about
113 rater behavior (14,31): What behaviors are assumed by the indices' estimations? What
114 behaviors in fact take place? Do the assumptions match the behaviors? The debaters rely on
115 theoretical arguments, mathematical derivation, fictitious examples, naturalistic comparisons,
116 and Monte Carlo simulation. A systematic observation of rater behavior is needed to inform
117 the debates over rater behavior.

118 This paper reports a controlled experiment that manipulated category, skew, and
119 difficulty, and observed raters' behavioral responses. The seven indices were tested against
120 the observed behavior. The findings also apply to the two equivalents of a_o , six equivalents of
121 S , two equivalents of π , and one equivalent of κ , covering 18 indices in total, all of which had
122 been analyzed mathematically by Zhao, Liu and Deng (14).

123 **Methods**

124 **Reconstructed Experiment with Golden Standard**

125 ***Reconstructed Experiment on Real Data (REORD)***

126 We conducted a $4 \times 8 \times 3$ between-subject controlled experiment with 4 subjects per
127 cell. Here the term “subject” refers to the unit of analysis of a study, such as a participating
128 patient in an experiment on the effectiveness of a new drug. A “subject” in this study,
129 however, was a rating session with 100 pairs of rating by two raters. As $4 \times 8 \times 3 \times 4 = 384$, this
130 study was based on 384 rating sessions, aka subjects. The three manipulated factors included
131 four levels of *category* ($C=2,4,6,8$), eight levels of *difficulty* (d_f ranges 0~1, 0 for the least
132 and 1 for the most difficult), and three levels of *skew* ($s_k=0.5$ for 50-50 distribution, 0.75 for
133 75-25 or 25-75 distribution, and 0.99 for 99-1 or 1-99 distribution), as summarized in Table
134 1.

135 [Insert Table 1 about here]

136 Over three hundred raters, registering 383 web names, from 53 Asian, European, and
137 North American cities judged online the lengths of bars, which served as the experimental
138 stimulus. A total of 22,290 items were rated, of which 19,900 were successfully paired,
139 producing 9,950 pairs of rating. Borrowing techniques from bootstrap (32,33), jackknife (34),
140 and Monte Carlo simulation (35), we sampled and resampled from the 9,950 pairs to
141 reconstruct the 384 rating sessions (36).

142 Thus, raters and rating were real, while rating sessions were reconstructed, making it a
143 *reconstructed experiment on real data* (REORD). The Appendix at the end of this manuscript

144 (Section II) provides further details and rationales of this REORD experiment.

145 ***Observed True Reliability (o_{ri}) and True Chance Agreement (o_{ac}) as Golden Standards***

146 The raters were instructed to judge the length of bars. The researchers determined the
147 bar lengths through programming, therefore know with certainty which rating decision was
148 right or wrong. As the lengths of the bars were set such that random guesses would occur
149 only between the longest and the second longest bars, the true chance agreement (o_{ac}) was
150 twice the wrong agreement (Eq. 3, Appx.), and true reliability (o_{ri}) was observed agreement
151 a_o minus o_{ac} (Eq. 5 of Appx.). Thus, o_{ri} served as the golden standard, namely the observed
152 estimand, against which the seven indices were evaluated, and o_{ac} served as the golden
153 standard for the seven chance estimators (37).

154 ***Five Independent Variables and Sixteen Dependent Variables***

155 Thus, this REORD experiment features three manipulated independent variables,
156 *category I*, skew (s_k) and difficulty (d_f) and 16 main dependent variables, which are the seven
157 indices' reliability and chance estimations plus the observed true reliability (o_{ri}) and true
158 chance agreement (o_{ca}). As the two main estimands, o_{ri} and o_{ca} sometimes also serve as
159 independent variables when assessing their impacts on the indices' estimations. Table 1,
160 Table 2 and the Appendix provide more details and rationales of variable calculations.

161 [Insert Table 2 about here]

162 ***Statistical Indicators – Directional R Squared (dr^2) and Mean of Errors (m_e)***

163 Reliability indices serve two functions. One is to evaluate measurement instruments
164 against each other, for which an index needs to accurately predict, meaning positively and
165 highly correlating with, true reliability. We use *directional r squared* ($dr^2=r\cdot|r|$) to gauge the
166 predictive accuracy of the seven indices and their chance estimators (Table 2 and Eq. 10 of
167 the Appendix). We preferred r^2 over r because r^2 has a clearer and more practical
168 interpretation, percent of the DV variance explained by the IV; r^2 is also more conservative as
169 $r^2\leq|r|$. We preferred dr^2 over r^2 because dr^2 indicates the direction of the relationship while r^2
170 does not.

171 The second function of the indices is to evaluate measurement instruments against
172 fixed benchmarks, such as 0.67 and 0.80, that some reliability authorities recommend
173 (19,30,38,39). For this function, an index needs to approximate true reliability. We use *mean*
174 *of errors*, m_e , which is the indices' deviations from the observed true reliability averaged
175 across the 384 rating sessions, to gauge the approximating accuracy of the seven indices,
176 denoted $m_e(r_i)$ in Table 2 and Eq. 8 of the Appendix. With the same reasoning, we also use m_e
177 to assess and compare the chance estimators of the indices, denoted $m_e(a_c)$ in Table 2 and Eq.
178 9 of the Appendix.

179 We adopted $dr^2>.8$ as the primary benchmark and $m_e<.02$ as the secondary benchmark
180 when evaluating the seven indices. Section V of the Appendix details the calculations of and
181 the rationales behind the benchmarks.

182 ***Functions of P Values and Statistical Pretests***

183 This study observes the tradition of reporting $p < \alpha$, where $\alpha = .05$, $.01$, or $.001$. We
184 however also strive to follow what have been advocated as a better statistical practice (40–
185 44):

- 186 1) avoiding the terms containing “significance, e.g., “statistical significance,” for
187 $p < \alpha$;
- 188 2) considering $p < \alpha$ as a prescreen threshold, passing which allows us to assess,
189 interpret, and compare effect size indicators, such as r^2 , dr^2 and m_e , with some
190 confidence;
- 191 3) using terms such as “statistical pretest” and “statistically acknowledged” where we
192 would have traditionally used “significance test” and “statistically significant;”
- 193 4) reserving the terms containing “significant” and “significance” exclusively for
194 effect sizes of practical or theoretical importance.

195 More of our views and practices regarding the functions of p values may be found in
196 our prior work (45–47).

197 **Results**

198 **Reliability Estimations Tested Against Observed Reliability**

199 Findings are summarized in Tables 3 through 6 and Figure 1 and discussed in three

200 sections. This section (II) reports the performance of the seven indices when predicting and
201 approximating the observed reliability. The next section (III) analyzes the impact of four
202 factors on the indices' performance. The following section (IV) discussed *offset* mechanism,
203 which is a key to understand the indices' complex behavior.

204 Overall, 2.86% of the raters' decisions fell on the short bars (1.11%, 1.93% and
205 5.53% respectively for four, six, and eight categories). As expected, there were fewer
206 agreements on short bars, averaging 0.45% (0.04%, 0.12%, and 1.18%). These agreements
207 showed no detectable effects on the main relations we investigate. The correlations between
208 the manipulated variables were practically zero, confirming orthogonality, which rules out
209 confounding or multicollinearity.

210 ***Predicting Reliability***

211 Percent agreement, a_o , the oldest and the most criticized index of interrater reliability,
212 did well predicting true reliability, showing $dr^2=.841$ (Line 3, Table 3). Of the seven indices
213 tested, a_o was the only one meeting the primary benchmark $dr^2>.8$ (Ineq. 11), outperforming
214 the second best, AC_1 ($dr^2=.721$), and the third best, S ($dr^2=.691$) by more than 10 points,
215 although the latter two met the tentative benchmark $dr^2>.67$.

216 [Insert Table 3 and Figure 1 about here]

217 The most respected three, π , κ and α , tied as the least accurate predictor, reporting
218 $dr^2=.312$, failing the tentative benchmark by margins. They also underperformed the next
219 worst, I_r , by 28.7 points ($dr^2=.599$).

220 The underperformances of the chance-adjusted indices, especially the popular π , κ and
221 α , were disappointing, considering that the whole mission of the indices was to outperform
222 a_o . The low r^2 means large predictive errors, suggesting that the three indices too often assign
223 lower scores to more reliable instruments, and attach higher scores to less reliable ratings.
224 They failed to differentiate reliable instruments from unreliable ones accurately and
225 consistently.

226 Figure 2 visualizes the performances and ranks the indices by their dr^2 scores. It is
227 noticed, again, that κ and α ranked among the lowest while percent agreement (a_o) ranked the
228 highest. Figure 2 also shows a strong and positive correlation between accuracy of predicting
229 chance agreement and accuracy of predicting interrater reliability ($dr^2=.9768$, $p<.001$),
230 supporting a design feature of this study, which is to analyze the indices' chance estimates for
231 the purpose of understanding the indices.

232 [Figure 2 About Here]

233 ***Approximating Reliability***

234 A .555 average reliability (o_{ri}) was observed (A3, Table 5). The seven indices'
235 estimation of reliability, however, ranged from .237 (π) to .726 (I_r), implying large

236 approximation errors. As expected, percent agreement (a_o) overestimated reliability, reporting
237 $e_m=.13$ (B6, Table 5) and $m_e=.13$ (A3, Table 4). The error, however, was below what's
238 allowed by the secondary benchmark, $m_e<.2$ (Ineq. 13 of the Appendix). So a_o was the only
239 index meeting both primary and secondary benchmarks.

240 [Insert Table 4 about here]

241 [Insert Table 5 about here]

242 Three other indices also met the $m_e<.2$ benchmark, of which two, AC_I ($m_e=.093$) and
243 S ($m_e=.096$). also outperformed a_o (Line 3 Table 4).

244 The trio, π , κ and α , again underperformed all others, reporting $m_e .323\sim.327$ (Line 8,
245 Table 5). The errors equaled one third of the 0~1 scale, and more than doubled the errors of
246 a_o ($m_e=.130$). I_r overestimated reliability across the board like a_o did (D6, Table 5), while κ , π
247 and α underestimated across the board -- 23.7%~24.1% estimated versus 55.5% observed
248 (Line 3, Table 5).

249 AC_I and S underestimated some sessions while overestimated other sessions (Line 6,
250 Table 5). Of AC_I and S , the under and over estimations offset each other to make the sizes
251 (absolute values) of e_m much smaller than that of m_e . Of the other five indices, e_m and m_e are
252 about equal in size (Line 6, Table 5 vs Line 3, Table 4).

253 In part because of the offsets, AC_I and S produced near-zero or very small e_m errors
254 (.001 and .044, respectively), much smaller than any of the other five indices did. By

255 contrast, κ , π and α again produced the largest errors, reporting e_m ranging from $-.318$ ~ $-.314$,
256 much worse than the next worst, I_r ($e_m=.171$, Line 6, Table 5).

257 ***Pi-Kappa-Alpha Synchrony***

258 As shown above, π , κ and α behaved like one index, despite the spirited debates on
259 which of them is the best (10,12,48–51). This pattern of π - κ - α *synchrony* persisted throughout
260 the data.

261 **Impacts of Four Factors**

262 The five viewpoints reviewed earlier discussed four factors behind reliability and/or
263 reliability estimations. Now that we have observed rater behavior, we examine the true
264 impacts of the four factors.

265 ***Conjecture Group 1: Chance Agreement Inflates a_o***

266 As said, a 13% chance agreement (o_{ac}) and a 55.5% reliability (o_{ri}) were observed,
267 while percent agreement (a_o) assumed 0% chance agreement and reported a 68.5% reliability,
268 which means a 13-point overestimation (Tables 4 and 5). Conjecture 1 and the century-old
269 beliefs were supported.

270 (1) Chance agreement exists.

271 (2) By completely overlooking chance agreement, a_o inflates the estimated reliability.

272 The data from this experiment, however, adds a third point:

273 (3) The chance agreement may not be as large as previously thought.

274 In this experiment, the chance agreement of a_o stayed below the .2 threshold, which
275 was a main factor that allowed the predictive accuracy (r^2) of a_o to stay above the .8
276 threshold. As a_o outperformed all six indices on the primary benchmark (r^2) and
277 outperformed four out of the six on the secondary benchmark (m_e), an argument could be
278 made that overestimating and misestimating chance agreement can be as counterproductive
279 as overlooking chance agreement.

280 ***Conjecture Group 2, Category Inflates S , I_r & AC_I***

281 As critics of S , I_r and AC_I would have predicted, $categoI(C)$ had large and negative
282 effects on chance estimations S_{ac} , I_{rac} and AC_{ac} , with dr^2 ranging -.863~- .661, ($p<.001$, Line
283 9, Table 3). Table 6 (K4~K7) shows more details, e.g., S_{ac} was 50% when $C=2$ but plunged to
284 12.5% when $C=8$. The decreases appeared large compared to the 13-point average o_{ac} .

285 [Insert Table 6 about here]

286 Negative effects on chance estimations contribute to positive effects on reliability
287 estimations, as shown in the dr^2 ranging .599 ~.721 ($p<.001$, Line 3, Table 3). S jumped from
288 40.2% when $C=2$ to 64.1% when $C=8$ ($C4\sim C7$, Table 6). The effect (difference) of 23.9
289 points is large compared with the 55.5-point average o_{ri} . In contrast, category effects on the
290 targets of estimations, o_{ri} and o_{ac} , were tiny. Coefficients dr^2 were respectively .003 ($p\geq.05$)
291 and -.019 ($p<.01$) (A4 and A9, Table 3, See Table 6, Lines 4~7, for more details).

292 These results support the classic theory that S and equivalents underestimate chance
293 agreement when categories exceed two, even when additional categories are largely empty.

294 The tables also show that I_r and AC_I relied on category in the same fashion that S did
295 and shared the same deficiency. The differences between the category effect on S , I_r or AC_I
296 estimation and the category effect on observed reliability all passed the $p < .001$ pretest. At the
297 meantime, category showed minimal effects ($dr^2 \approx .001$, $p \geq .05$) on π , κ and α , as their authors
298 intended (Line 4, Table 3).

299 ***Conjecture Group 3: Skew Depresses κ , π & α***

300 As critics of κ , π & α would have predicted, skew had substantial and positive effects
301 on chance estimators κ_{ac} , π_{ac} & α_{ac} , with dr^2 ranging .434~.437 ($p < .001$, Line 10, Table 3).
302 Table 6 (Lines 8~10) shows more details, e.g., κ_{ac} was 50% when distribution was 50&50,
303 but rose to 67.6% when distribution changed to 1&99.

304 The positive effects on chance estimates led to negative effects on reliability
305 estimates. Skew effects on the three indices were all negative, with dr^2 ranging $-.293 \sim -.292$
306 ($p < .001$, Line 5, Table 3). When distribution changed from completely even to extremely
307 skewed, the trio's chance agreement estimates increased from about .5 to about .68, and in
308 parallel their reliability estimates decreased from about .37 to about .04, a drop of over 89%
309 (Lines 8~10, Table 6). While mathematical analyses of prior studies had predicted a drop
310 (14,26,52), the empirical evidence of this study showed the drastic magnitude of the drop.

311 In contrast to the large effects on the index estimators, skew showed minimal effect
312 on the observed estimands, o_{ri} and o_{ac} ($p \geq .05$ for both dr^2 , A5 & A10, Table 3), supporting the
313 argument that chance estimates and reliability indices should not rely on skew. Each
314 difference between the skew effect on π , κ or α estimation and the category effect on the
315 observed estimand passes the $p < .001$ pretest.

316 In another contrast, skew showed practically zero effects on S , I_r or their chance
317 estimates, and a small negative effect on AC_{ac} ($dr^2 = -.039$, $p < .001$, Lines 5 & 10, Table 3). So
318 I_r avoided the skew effect as its authors intended, while AC_I reversed the effect as its author
319 intended, although the reversed effect was small. A long-suspected pattern was confirmed
320 empirically -- κ , π & α were dependent on skew while S , I_r & AC_I were dependent on
321 category.

322 ***Conjecture Group 4: Indices Overlook Task Difficulty***

323 *Difficulty* showed a substantial and positive effect on o_{ac} ($dr^2 = .585$, $p < .001$, A11,
324 Table 3), and a large and negative effect on o_{ri} ($dr^2 = -.774$, $p < .001$, A6). A change from
325 extremely easy to extremely difficult decreased o_{ri} by over 68 percentage points and
326 increased o_{ac} by nearly 36 points (Columns A and I, Table 6). These effects appear large
327 compared with 13-point average o_{ac} and 55.5-point average o_{ri} , suggesting that chance
328 estimates and reliability indices should rely on difficulty.

329 In contrast, difficulty had minimal effects on S_{ac} , I_{rac} and AC_{ac} ($dr^2=.000\sim.009$, $p\geq.05$,
330 Table 3) and *negative* effects on κ_{ac} , π_{ac} & α_{ac} ($dr^2=-.123$ or $-.125$, $p<.001$, Table 3; c.f.
331 Columns I & N~P, Lines 11~18, Table 6), implying that the indices either failed to rely on
332 difficulty or relied on its opposite, *easiness*, to estimate chance agreement. Each difference
333 between the difficulty effect on chance estimation and the difficulty effect on observed
334 chance agreement was statistically acknowledged at $p<.001$.

335 Difficulty showed weaker effects on the six chance-adjusted indices ($dr^2=-.566\sim-.389$,
336 Line 6, Table 3) than on the estimation target o_{ri} ($dr^2=-.774$). Each difference between the
337 difficulty effect on reliability estimation and the difficulty effect on observed reliability was
338 statistically acknowledged at $p<.001$.

339 By contrast, a_o , showed a strong and negative correlation ($dr^2=-.778$, B6, Tables 3)
340 with difficulty. The correlation was as strong as the correlation between o_{ri} and difficulty
341 ($dr^2=-.774$, A6), suggesting the negative correlations between the chance-adjusted indices
342 and difficulty ($dr^2=-.566\sim-.389$) are likely due to a_o embedded in the indices.

343 Based on derivation and simulation, Gwet concluded that the indices before AC_I had
344 not handled difficulty properly, and AC_I handled it better, at least than κ (53–55). The above
345 findings support both claims. The near zero correlation between AC_{ac} and difficulty
346 ($dr^2=.009$, $p\geq.05$, E11, Table 3), however, suggests that AC_I still does not handle difficulty
347 well.

348 ***Conjecture Group: Indices Assume Intentional and Maximum Random Rating***

349 The precision evidence for the behavioral assumptions behind the statistical indices
350 comes from mathematical analysis. A 2013 study provides detailed scenarios of rater
351 behavior assumed by each of the 22 indices analyzed (14). Readers are invited to derive
352 mathematical formulas from the behavioral scenarios. If a reader-derived formula matches
353 the formula for the corresponding index, then the reader may conclude that the
354 corresponding index indeed assumes the behavioral pattern spelt out in the scenario. If, for
355 example, a formula derived from the Kappa Scenario provided by the 2013 study matches the
356 formula for Cohen's κ published in 1960 (2), it would confirm that κ indeed assumes the rater
357 behavior depicted in the 2013 Kappa Scenario. Such exercises by readers have shown them
358 that chance-adjusted indices all assume that raters regularly conduct *intentional and*
359 *maximum random rating*.

360 This study provided corroborating empirical evidence. The indices' chance estimates
361 were poorly correlated with their estimands, the observed chance agreements (Table 3, Line
362 8). The observed chance agreement (o_{ac}) explained less than 8% of the variance in each of the
363 category-based indices' chance estimates, S_{ac} (2.1%), I_{rac} (2.1%), and AC_{ac} (7.5%). Although
364 the correlations were stronger for the skew-based indices' chance estimates, π_{ac} (-15.1%), κ_{ac}
365 (-15.2%), and α_{ac} (-15.1%), the dr^2 coefficients were all negative, suggesting that the three
366 indices tended to give higher estimates when the true chance agreements were lower, and

367 give lower estimates when the true chance agreements were higher. Clearly, the index-
368 estimated random ratings were not the raters' random ratings observed in this study. This
369 finding supports the argument that the chance-adjusted indices assume intentional and
370 maximum random rating while typical raters conduct involuntary and task-dependent random
371 rating. The mismatch between the assumptions and the observations explains the negligible
372 or negative correlations between the estimates and the estimands.

373 More corroborating evidence for the maximum-random assumption came from the
374 large overestimation of chance agreement by the six chance-adjusted indices, as shown at
375 Line 12 of Table 5 and the right half of Table 6, summarized in Line 19.

376 The more situational and detailed evidence of the behavioral assumptions come from
377 the influences of the four factors and the offset and aggravation behaviors of the indices,
378 which are discussed below.

379 ***Summarizing the Impact of Four Factors***

380 Each index of interrater reliability implied one or more misassumptions about chance
381 agreement. a_o Overlooked chance agreement. S , I_r and AC_l inappropriately relied on
382 category. π , κ And α inappropriately relied on skew. While difficulty had a strong and
383 positive effect on chance agreement, all chance adjusted indices failed to rely on difficulty. π ,
384 κ and α even relied on its opposite, easiness. The misassumptions, including missed,
385 mistaken, and contra assumptions, impeded estimation. π , κ and α fared worse in part because

386 they entailed more and more devastating misassumptions, some of which had been mistaken
387 as signs of sophistications.

388 Recall that the main mission of chance adjusted indices is to remove chance
389 agreement in order to improve on percent agreement. When they mishandled the factors
390 affecting chance agreement, they misestimated chance agreement, thereby misestimated
391 reliability. Misassumptions about the four factors are keys to understanding the indices'
392 underperformance.

393 To understand more, we discuss below the *offsetting* mechanism, which interacts with
394 the assumptions and misassumptions of the indices to define the indices' behavior.

395 **Offsets in Reliability Estimation**

396 Puzzles may arise if one peruses Tables 3 through 6, five of which discussed below.

397 **Puzzle 1.** Each chance-adjusted index relied on a wrong factor, skew or category, to
398 estimate chance agreement; none of them relied on a right factor, difficulty. How come some
399 approximated chance agreement far better than others (Line 12 of Table 5 and Line 7 of
400 Table 4)?

401 **Puzzle 2.** Chance estimators barely measured the observed chance agreement o_{ac} , or
402 even measured anti o_{ac} (C8~H8 of Table 3). How come the reliability estimations were all
403 positively and sometimes substantially correlated with the observed reliability (C3~H3)?

404 **Puzzle 3.** Assuming a negative relation between chance agreement and reliability, one

405 might expect that an over estimation of chance agreement leads to an under estimation of
406 reliability. How come S overestimated chance agreement by 100% ($o_{ac} = .130$ compared to
407 $S_{ac} = .260$, Line 9, Table 5) while at the same time approximated reliability almost perfectly
408 ($S = .556$, compared to $o_{ri} = .555$, Line 3, Table 5)?

409 **Puzzle 4.** Continued from Puzzle 3, how come AC_I overestimated chance agreement
410 ($e_m = .044$, Line 12, Table 5) while also overestimated reliability ($e_m = .044$, Line 6, Table 5)?

411 More generally, how come across-the-board overestimations of chance agreement did
412 not translate into across-the-board underestimations of reliability (Line 12 vs Line 6, Table
413 5)?

414 **Puzzle 5.** Continued from Puzzles 3 & 4, how come I_r overestimated chance
415 agreement more than AC_I did ($I_{rac} = .131$ vs $AC_{ac} = .044$, Line 12, Table 5), while also
416 overestimated reliability more than AC_I did ($I_r = .171$ vs $AC_I = .044$, Line 6, Table 5)?

417 The puzzles can be explained in part by *offsets*, including *partial offset*, *over offset*,
418 and *counter offset (aggravation)* built into the reliability formulas, some of which discussed
419 below.

420 ***Category offset, skew aggravation, and skew offset***

421 To understand Puzzle 1, first recall that, under intentional-and-maximum-random
422 assumption, chance-adjusted indices tend to overestimate chance agreement
423 (9,14,29,38,39,56–58). In this experiment, the overestimations ranged from 4.4 percentage

424 points by AC_I to 44.5 points by Scott's π , all statistically acknowledged ($p < .001$, Line 12,
425 Table 5).

426 To explain Puzzle 1, we note that the category-based indices assume that larger
427 number of categories *decreases* chance agreement (C9~E9, Table 3), which *offset* the general
428 overestimation. The skew-based indices assume that higher skew *increases* chance agreement
429 (F10~H10), which *aggravated* the general overestimation. AC_I assume both, that category
430 and skew both decrease chance agreement (E10), thereby *offset* the overestimation even more
431 than the other two category-based indices.

432 To illustrate the point, we follow the textbook tradition of starting from *ground zero*,
433 which is the condition of two raters, two categories, and 50&50% distribution. Here, and only
434 here, all major indices gave about the same estimates, $a_c \approx 0.5$ (K2~P2, Table 6). Under
435 intentional-and-maximum-random assumption, two raters draw from marbles, half with one
436 color and half another color; they rate randomly if the colors match, and honestly if mismatch
437 (9,14,29,38,39). Task difficulty is not a factor in this view of rater behavior.

438 In actual rating, however, $a_c = 0.5$ could occur only if the task is extremely difficult. In
439 our experiment, even the most difficult ($d_f = 1$ for 1-pixel difference) condition did not reach
440 that theoretical maximum, reporting an $o_{ac} = .38$ (I18, Table 6). The less difficult sessions
441 reported significantly smaller o_{ac} , averaging 0.13 across all levels of difficulty. This means a
442 37-point initial overestimation at the ground zero by each chance-adjusted index

443 ($e_m=.5-.13=.37$).

444 When category increased from ground zero, S_{ac} , Ir_{ac} and AC_{ac} decreased quickly
445 under the *category assumption* (Columns K~M, Row 4~7, Table 6). While the assumption
446 was unjustified given the small change in o_{ac} (I4~I7), the decrease partially offset the 37-
447 point overestimation, making S_{ac} , Ir_{ac} and AC_{ac} less inaccurate. By contrast, κ_{ac} , π_{ac} & α_{ac}
448 rejected the category assumption to remain unchanged (Columns N~P), hence did not benefit
449 from the partial offset. Thus, S_{ac} , Ir_{ac} & AC_{ac} became less inaccurate than κ_{ac} , π_{ac} & α_{ac} .

450 Now return to ground zero, then increase skew. Under the skew assumption, κ_{ac} , π_{ac} &
451 α_{ac} increased with skew (Columns N~P, Row 8~10, Table 6). While the assumption was
452 unjustified given the small change in o_{ac} (I8~I10), the increase further aggravated the 37-
453 point overestimation, making κ_{ac} , π_{ac} & α_{ac} even more inaccurate. By contrast, S_{ac} and Ir_{ac}
454 rejected the skew assumption to remain unchanged (K~L, 8~10), hence did not suffer from
455 the aggravation. Thus, κ_{ac} , π_{ac} & α_{ac} became even more inaccurate than S_{ac} & Ir_{ac} .

456 Rather than accepting or rejecting the skew assumption, AC_{ac} reversed it, by assuming
457 that skew reduced a_c (M8~M10). While the assumption also mismatched the observed skew
458 effects (I8~I10), the decrease further reduced the once 37-point overestimation. Here two
459 unjustified assumptions, *category* and *reversed skew*, joined hands to partially offset another
460 unjustified assumption, *intentional and maximum random*. Thus, AC_{ac} became even less
461 inaccurate than S_{ac} & Ir_{ac} , hence the least inaccurate of the six. As the effect of intentional-

462 and-maximum-random assumption was stronger than the other two effects combined, a net
463 effect was that even *ACac* still overestimated chance agreement.

464 There were other under-offsets, over-offsets, and counter-offsets, i.e., aggravations,
465 some of which discussed below. Behind multifarious offsets were multifarious assumptions
466 about rater behaviors, which fought or allied with each other or stayed neutral to produce the
467 multifarious outcomes. Two wrongs sometimes made one right, sometimes half right, and
468 often three, four, or more wrongs.

469 *Chance-removal offset*

470 To understand Puzzle 2, we first recall that, assuming intentional and maximum
471 random rating, index designers want to remove maximum amount of chance agreement from
472 all considerations, which requires to remove a_c not only from percent agreement (a_o), but also
473 from the realm of consideration (9,14,23,24,29,38,39). Accordingly, a_c is subtracted twice in
474 Eq. 1, first from a_o in the numerator, and second from 1 in the denominator, which represents
475 100% of the realm of consideration. Two offsets occurred as a result. First, a_c offsets a_o in the
476 numerator. Second, a_c in the denominator offsets its own impact in the numerator. As the
477 self-offsets weaken a_c 's effects, a_o dominates Eq. 1, the indices' estimation of reliability. That
478 explains Puzzle 2: the minimal or negative a_c-o_{ac} correlations exerted weaker effects than the
479 strong and positive a_o-o_{ri} correlation.

480 The weaker effects still hinder. The chance estimators not only failed to fulfill their

481 prescribed mission of improving on percent agreement, but the estimators worked against the
 482 mission. Consequently, all six indices underperformed percent agreement when predicting
 483 observed true chance agreement. Ironically, it was the supposedly “most primitive” and
 484 “flawed” percent agreement (a_o) that worked inside the indices to keep them from performing
 485 and looking even worse (2 p38,12 p80).

486 The offsets also help to explain Puzzle 3. While S overestimated chance agreement by
 487 an averaged 13.1 points (Line 12, Table 5), the chance-removal offset helped to bring down
 488 the scalar error of reliability estimation to 9.6 points (Line 3, Table 4). This across-session
 489 error contains over- and under-estimations of individual sessions, which offset each other in
 490 averaging to reduce the vector error to near zero ($e_m=.001$, Line 6, Table 5. See also the
 491 discussion of aggregation bias earlier).

492 By setting estimated reliability (r_i in Eq. 1) equal to observed reliability (o_{ri} in Eq. 5
 493 of Appendix), $r_i=o_{ri}$, we derive a threshold (t_h) for a_c , which is Eq. 2:

$$t_h = \frac{o_{ac}}{1 - o_{ri}} \quad 0 \leq o_{ac} \leq t_h \leq \infty \quad (2)$$

494 For any rating session, an index accurately estimates reliability when $a_c=t_h$,
 495 underestimates when $a_c>t_h$, and overestimates when $a_c<t_h$. Therefore, when $o_{ac}<a_c<t_h$, the
 496 index overestimates both the chance agreement and the reliability, explaining Puzzle 4.

497 Across the 384 sessions, average t_h would be .292 if we plug o_{ac} (.13) and o_{ri} (.555) into Eq.
 498 2. As Table 5 shows, of the six chance-adjusted indices, the three (κ , π , α) reporting $a_c>.292$

499 (Line 9) also underestimated reliability (Line 6), and the three (S , I_r , AC_I) reporting $a_c < .292$
500 also overestimated reliability. At the same time, all six overestimated chance agreement (Line
501 12). Due to the chance-removal offset, it is possible and possibly common for some category-
502 based indices to overestimate both chance agreement and reliability.

503 A previously undocumented paradox emerges from this analysis (Eq. 1 and Eq. 2). An
504 index estimates reliability accurately ($r_i = o_{ri}$) *only* when it overestimates chance agreement
505 ($a_c > o_{ac}$), an index that estimates chance agreement accurately ($a_c = o_{ac}$) inevitably
506 underestimates reliability ($r_i < o_{ri}$), except in the extreme and impractical situation when
507 $r_i = o_{ri} = 0$. The paradox, applicable for all known chance-adjusted indices, is rooted in the
508 chance-removal offset imposed by Eq. 1, which traces back to the intentional and maximum
509 random assumption (14,23,24,26).

510 ***Square-root over offset***

511 To understand Puzzle 5, recall that Perreault and Leigh's I_r adopts the chance
512 estimator of S , $I_{rac} = S_{ac}$, and takes the square root of S as the reliability estimation (7). $S \leq I_r$, as
513 $I_r = S^{1/2}$ for $1 \geq S \geq 0$ and $I_r = 0$ for $-1 \geq S < 0$. When chance agreement is overestimated, the square
514 root operation constitutes an additional offset (14). Due to the category-based over-offset of
515 S , I_r overestimates chance agreement more than AC_I ; at the meantime, due to the square root
516 over-offset of I_r , I_r overestimates reliability more than AC_I . The two offsets explain Puzzle 5.

517 A rating session in this experiment simulates a study. In practice, errors do not offset

518 across studies, e.g., one study's overestimation of Disease A does not offset another study's
519 underestimation of Disease B. We should not overemphasize the near-zero aggregated error
520 by S shown in e_m or overlook the sizable individual errors by S shown in m_e .

521 **Discussion**

522 **Main Findings**

523 Of the seven indices, percent agreement (a_o) stood out as the most accurate predictor
524 of reliability ($dr^2=.841$, Table 3) and the third most accurate approximator ($m_e=.130$, Table
525 4). AC_I , the newest and the least known, was the second-best predictor ($dr^2=.721$) and the
526 best approximator ($m_e=.093$). S ranked behind AC_I for both functions ($dr^2=.691$, $m_e=.096$).

527 The most respected, the most imposed, and the most applied indices, π , κ and α ,
528 ranked the last for both functions ($dr^2=.312$, $m_e=.323\sim.327$).

529 The indices' underperformances appeared attributable to mismatches between the
530 assumed and observed rater behaviors, and multifarious offsets and aggravations between the
531 misassumptions. Percent agreement assumed zero random rating, leading to the 13-point
532 overestimation of reliability. The other six indices assumed intentional and maximum random
533 rating, leading to a 37-point initial overestimation of chance agreement at "ground zero" for
534 interrater reliability (Line 3, Table 6).

535 Away from ground zero, S , I_r and AC_I assumed larger number of categories produced
536 less chance agreement, which offset the initial overestimation, while π , κ and α assumed

537 skewer distributions produced more chance agreement, which aggravated the overestimation.
538 The opportune offsets and the austere aggravations explain the smaller approximation errors
539 by the category-based indices than by the skew-based indices. Contrary to the assumptions,
540 neither rating category nor distribution skew showed meaningful effects on the observed true
541 chance agreement.

542 Difficulty exhibited a substantial and positive effects on chance agreement ($dr^2=.585$,
543 $p<.001$, Table 3), while S , I_r , and AC_I did not rely on difficulty to estimate chance agreement
544 ($dr^2=.000\sim.009$, $p\geq.05$). Failing to rely on difficulty further explains the three indices'
545 underperformance in prediction. Moreover, π , κ & α relied on the opposite, easiness, to
546 estimate chance agreement ($dr^2 =-.125\sim-.123$, $p<.001$), which contributed another part to π , κ
547 & α 's worse performance than S , I_r , and AC_I .

548 **What Did the Indices Indicate?**

549 An index indicates a certain concept. What did the seven indices indicate? Did they
550 indicate what they purport to indicate?

551 Percent agreement a_o was the only index meeting the primary benchmark ($dr^2>.8$),
552 thereby also meeting the competitive benchmark. By overlooking chance agreements, a_o
553 overestimated reliability by 13 percentage points ($e_m=m_e=.130$, Tables 4 & 5). The error,
554 however, was within the margin allowed by the secondary benchmark ($m_e<.2$). The
555 overestimation appeared across the board, as shown in Columns A and B (Lines 4 through

556 18) of Table 6, which implies that researchers and reviewers may manage a_o 's deficiency by
557 discounting a certain amount, such as 15 points, treating $a_o-0.15$ as a crude estimation of
558 reliability. Overall, in this experiment percent agreement behaved as a good predictor and a
559 13-point over-approximator of interrater reliability.

560 The other six indices set out to outperform a_o by removing estimated chance
561 agreement a_c . Unfortunately, their a_c estimations failed to accurately estimate true chance
562 agreement o_{ac} . S_{ac} , Ir_{ac} , and AC_{ac} were slightly influenced by o_{ac} ($dr^2=.021\sim.075$, $p<.01$ or
563 $p<.001$, Table 3). They were instead strongly and negatively influenced by category
564 ($dr^2=-.863\sim-.661$, $p<.001$), suggesting they indicated fewness of category more than they
565 indicated chance agreement. The other three chance estimators, π_{ac} , κ_{ac} & α_{ac} , predicted far
566 less accurately. They indicated mostly skew ($dr^2=.434\sim.437$) and, to a lesser extent, easiness,
567 the opposite of o_{ac} (Lines 8-10, Columns F-H, Table 3).

568 When Eq. 1 was used to remove a_c , a_o offset some impact of a_c , which also self-offset
569 some. The offsets reduced the category and skew effects and kept the index- o_{ri} correlations
570 positive (Line 3-5, Table 3). But still, a_c , the unique core of each index, all impeded the
571 reliability estimation. S_{ac} , Ir_{ac} and AC_{ac} impeded less than π_{ac} , κ_{ac} , & α_{ac} did, allowing S , I_r
572 and AC_I to predict reliability better than π , κ , & α did (Line 3, Table 3). But the reduced
573 impediments were still impediments. Consequently, none of the chance-adjusted indices had
574 a good chance of outperforming a_o when predicting reliability. Two indices, AC_I ($m_e=.093$)

575 and S ($m_e=.096$), did outperform a_o ($m_e=.13$) for approximation, which was due more to
576 opportune offsets between misassumptions, and less to removing chance agreements (Line 3,
577 Table 4).

578 At the end, no chance-adjusted index passed the primary benchmark $dr^2>0.8$. Two,
579 AC_I (.721) and S (.691), passed the threshold $dr^2>0.67$ for tentative acceptance (Table 3).
580 Being the best approximator, AC_I ($m_e=.093$) was the one meeting the competitive benchmark.
581 AC_I and S were also two of the four indices meeting the secondary benchmark, $m_e<.2$ (Line 3,
582 Table 3).

583 Category exerted some effects on AC_I ($dr^2=.123$) and S ($dr^2=.175$). Fortunately for the
584 two indices, the category effects were much smaller than the estimand effects of o_{ri} ($dr^2=.721$
585 & .691). The two indices underestimated reliability when $C=2$, and overestimated when $C\geq 4$
586 (Columns A, C and E, Lines 4~7, Table 6). Overall, AC_I and S were acceptable predictors of
587 interrater reliability, and under- or over-approximators when category was respectively under
588 or over 3.

589 I_r ($dr^2=.599$, $m_e=.18$) failed the tentative benchmark for prediction but satisfied the
590 secondary benchmark for proximity. It overestimated reliability across the board. Overall, I_r
591 was a poor predictor and an 18-point over-approximator of interrater reliability. I_r 's
592 overestimation was worse when the number of categories was increased.

593 The performances of π , κ and α belong to another class. The trio's estimation-
594 estimand correlations ($dr^2=.312$) were far below the primary benchmark of $dr^2>.8$ or the
595 tentative benchmark of $dr^2>.67$; and their approximation errors ($m_e=.323\sim.327$) were far
596 above the secondary benchmark $m_e<.2$. Furthermore, evenness (1-skew) exerted nearly as
597 large effects on the trio ($dr^2=.292\sim.293$, Line 5) as their estimand o_{ri} did ($dr^2=.312$),
598 suggesting that the trio indicated distribution evenness nearly as much as they indicated
599 interrater reliability. More even distributions raised π , κ and α nearly as effectively as higher
600 reliability did, even though skew or evenness showed no effect on observed reliability or
601 chance agreement.

602 Overall, π , κ & α were crude predictors of reliability and evenness, and 31-point
603 under-approximators of reliability. They were crude because they showed large errors when
604 predicting reliability ($dr^2=.312$) or evenness ($dr^2=.292\sim.293$).

605 While $dr^2 (.292\sim.293)$ were too low to make π , κ & α precise indicators of evenness
606 or skew, they were too high to allow the trio to be pure indicators of reliability. The
607 correlation can be even more disconcerting if one considers its impact on the creation and
608 selection of scientific knowledge. Reviewers and researchers use the trio to screen
609 measurements and manuscripts, while trio systematically favor more even distributions,
610 making the world appear flatter. It would be a collective version of the conservative bias,

611 except this one permeates scientific knowledge (59,60). By contrast, a_o showed none of this
612 disparaging deficiency ($dr^2=.000$).

613 **Conclusion**

614 Like most controlled experiments, this study had limited external validity. The raters
615 made visual judgments, which did not represent all tasks. The categories stopped at eight.
616 The short-bar categories were largely empty by design. Each session had only two raters. The
617 list could go on. To avoid unwarranted generalization, we used past tense to describe the
618 indices' behaviors and their impact.

619 Our findings, however, have been speculated or predicted by the theoretical analyses,
620 mathematical derivations and Monte Carlo simulations (14,29,64,65,53,55–58,61–63). These
621 studies used no actual measures, specific tasks, human raters, or other specifics that may limit
622 external validity. What some other studies lack in internal validity, this study provides. The
623 validity of our collective knowledge is significantly strengthened by adding empirical studies
624 based on observing rater behavior.

625 The indices were advertised to be “standard” and “global” for “general purpose”
626 (12,14,66,67). Now that some reigning indices did not perform as advertised against one set
627 of observed behavior, it is good evidence that indices are not general or global or standard.
628 The burden is not on doubters to prove that the indices always fail, but on defenders to
629 demonstrate that the indices perform, at least sometimes.

630 Despite the lack of empirical evidence in support of the reigning indices, the spiral of
631 inertia may continue, forcing some and enticing others to work with the indices (26,52). In
632 that event, the interpretation of π , κ and α may warrant more caution, and the application of
633 a_o and AC_I may deserve more credence, to the extent that findings of this experiment will be
634 replicated.

635 **Future Research**

636 *Replication studies.* More controlled experiments are called for to falsify or qualify
637 the findings of and the theories behind this experiment, and to test the other reliability indices
638 against their estimands (66,68,69).

639 *New Indices.* New indices may be needed. Index designers may be more cautious
640 about the assumptions that raters conduct intentional and maximum chance rating, or their
641 chance rating is determined by skew or category. More thoughts may be given to the
642 possibility that raters conduct instead involuntary and task-dependent random rating, and
643 more weights given to task difficulty. The index designers are encouraged to assess and
644 adjust their ideas and indices against behavioral data, including the data from this experiment,
645 which will be made public upon publication of this manuscript.

646 *REORD and Behavior-based statistical methods.* Mathematical statistics use a
647 system of axioms and theorems to build tools for analyzing behavioral data. The REORD
648 (reconstructed experiment on real data) methodology reverses the logic, using observed

649 behavior to inform statistical methods. The application might not be limited to interrater
650 reliability. REORD, for example, may open a new front for the studies of sensitivity and
651 specificity measures, two practical tools often used in medical and health research. REORD
652 may also help to investigate the empirical relationship between reliability and validity, two of
653 the most fundamental concepts of scientific enquiry.

654 ***Rater expectations of prevalence or skew.*** The researchers in this REORD
655 experiment told the raters nothing about the prevalence or the skew of the long and short
656 bars. As prevalence and skew were programmed to vary randomly between trials and
657 between rating sessions, the researchers themselves did not know about the prevalence or
658 skew until data analysis, and the raters could not have guessed accurately. This design feature
659 was chosen because it resembled one type of research condition, under which raters don't
660 know what to expect, therefore they don't expect.

661 For some tasks, however, raters do expect about prevalence and skew, due to their
662 prior experience with the same tasks or their prior exposure to second-hand information. A
663 follow-up study may investigate the impact of such expectations on raters' rating or the
664 indices of reliability, sensitivity, and specificity.

665 ***Human vs machine raters.*** Expectations about distribution, prevalence, and skew
666 can be programmed into artificial intelligence (AI) to aid automated diagnoses, judgements,
667 scorings, evaluations, ratings, and other decisions by machines. Unlike human decisions and

668 human expectations that are often vague and varying, machine decisions and machine
669 expectations can be programmed to be super clear and super consistent (70,71). Topics of
670 human-machine reliability and inter-machine reliability versus inter-human reliability could
671 be fruitful and fascinating for research using REORD, and so could topics of sensitivity,
672 specificity, and validity with human and/or machine raters.

673

674 **Declarations**

675 **Ethics approval and consent to participate**

676 The survey study received ethical approval under the ethics procedures of University of
677 Macau Panel on Research Ethics (reference SSHRE22-APP016-FSS). Written consent for the
678 survey was also taken.

679 All methods were carried out in accordance with relevant guidelines and regulations.

680 Informed consent was obtained from all subjects/participants and/or their legal guardian(s).

681 **Consent for publication**

682 Not applicable.

683 **Availability of data and materials**

684 The datasets used and/or analyzed during the current study are available from the
685 corresponding author on reasonable request.

686 **Competing interests**

687 The authors declare that they have no competing interests.

688 **Funding**

689 This research is supported in part by grants of University of Macau, including CRG2021-
690 00002-ICI, ICI-RTO-0010-2021, CPG2021-00028-FSS and SRG2018-00143-FSS, ZXS PI;
691 Macau Higher Education Fund, HSS-UMAC-2020-02, ZXS PI; Jiangxi 2K Initiative through
692 Jiangxi Normal University School of Journalism and Communication, 2018-08-10, Zhao PI.

693 **Acknowledgements**

694 The authors gratefully acknowledge the contributions of Hui Huang and Chi Yang to the
695 execution of the reconstructed experiment.

696 **Authors' contributions**

697 XZ designed the study, supervised the construction of the experimental site, organized the
698 data collection, conducted the data analysis, and drafted the manuscript. GCF provided
699 feedback for the research design, assisted with data analysis, and provided comments. SHA
700 and LPL provided input into the manuscript writing. All authors read and approved the final
701 manuscript.

702

703

704	Appendix to “Interrater reliability estimators	
705	tested against true interrater reliability”	
706	I:	Five Concepts and Four Viewpoints about Interrater Reliability
707	I.1.	Five Concepts
708		I.1.1. Interrater Reliability (r_i)
709		I.1.2. Chance Agreement (a_c)
710		I.1.3. Category (C)
711		I.1.4. Distribution Skew (s_k)
712		I.1.5. Difficulty (d_f)
713	I.2.	Five Viewpoints
714		I.2.1. Chance agreement inflates a_o
715		I.2.2. Rating category inflates S , I_r , and AC_I
716		I.2.3. Distribution skew deflates π , κ & α
717		I.2.4. Reliability indices overlook task difficulty
718		I.2.5. Indices assume intentional and maximum random rating
719	II:	Reconstructed Experiment with Golden-Standard Task
720	II.1.	Manipulating Category (C)
721	II.2.	Manipulating Difficulty (df)
722	II.3.	Creating One-way Golden Standard
723	II.4.	Pairing Rater Responses
724	II.5.	Manipulating Skew (s_k)
725	II.6.	Reconstructing Rating Sessions
726	II.7.	Reconstructed Experiment in Summary
727	III:	Variable Measurements and Calculations
728	III.1.	Calculating Chance Agreement (o_{ac})
729	III.2.	Alternative Calculation of Observed Chance Agreement (o_{ac})
730	III.3.	Calculating Observed Reliability (o_{ri})
731	IV:	Statistical Indicators
732	IV.1.	Approximating and predictive functions of reliability indices
733	IV.2.	Proximity Measure I -- Error of Mean (e_m)
734	IV.3.	Proximity Measure II -- Mean of Errors (m_e)
735	IV.4.	Predictive Accuracy and Share of Influence -- Directional r^2 (dr^2)
736	IV.5.	Regression vs ANOVA

737	V:	Benchmarks and Thresholds
738		V.1. Ideal index outperforms all others
739		V.2. Reliability over chance agreement
740		V.3. Prediction (dr^2) over approximation (m_e & e_m)
741		V.4. m_e over e_m
742		V.5. Primary Requirement
743		V.6. Secondary Requirement
744		V.7. Tentative Requirement
745		V.8. Competitive requirement

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765 **I. Five Concepts and Five Viewpoints about Interrater Reliability**

766 Section I discusses five concepts and five viewpoints about interrater reliability to
767 supplement the literature review and design sections of the manuscript.

768 **I.1. Five Concepts**

769 Five of the fundamental concepts, interrater reliability (r_i), chance agreement (a_c),
770 rating categories (C), distribution skew (s_k), and task difficulty (d_f), is explicated below (72).
771 Indicators of r_i and a_c were measured as dependent variables in this experiment, while C , s_k
772 and d_f were manipulated as the three independent variables.

773 **I.1.1. Interrater Reliability (r_i).** *Interrater reliability* (r_i) refers to the true agreement
774 between raters, aka coders, engaged in systematic and task-driven rather than random rating,
775 aka coding. Indices of interrater reliability are meant to estimate this true agreement. As
776 chance agreement (a_c), defined below, is believed to inflate reliability estimate, many indices
777 attempts to estimate and remove a_c (7,53,55,73,74). All major indices, including the six
778 examined in this study other than %-agreement a_o , share Eq. 1 to remove a_c and estimate r_i .

$$r_i = \frac{a_o - a_c}{1 - a_c} \quad (1)$$

779 A main objective of this study is to assess the seven indices of interrater reliability
780 against observed true reliability (o_{ri}). The eight measures also serve as dependent variables,
781 on which the effects of category, skew, and difficulty are assessed and compared.

782 **I.1.2. Chance Agreement (a_c).** This study also measured seven chance agreement (a_c)
783 variables, one chance estimate for each of the six chance-adjusted indices plus observed true
784 chance agreement (o_{ac}). There was an implied eighth chance indicator, by percent agreement
785 (a_o), which is, by definition, a constant at zero.

786 *Chance agreement (a_c)* refers to the agreement produced by random rather than
787 systematic and task-driven rating. Five indices invented their own chance estimators while I_r
788 adopted the estimator from S .

789 In addition to comparing the indices with observed reliability, it is important to also
790 compare the indices' chance estimates with observed chance agreement. In Equation 1,
791 subtraction of a_c in the nominator decreases r_i , while the subtraction in the denominator
792 increases r_i . The varying offsetting obscures the differences between indices (14). Since the
793 main or only difference between many indices is in chance estimators (a_c), comparing a_c with
794 its estimation target (o_{ac}) may tell us more about the inside mechanism at the core of the
795 indices.

796 The seven chance estimate measures also serve as dependent variables, on which the
797 effects of category, skew, and difficulty are assessed and compared.

798 **I.1.3. Category (C).** *Category (C)* was defined as the number of choices available to a
799 rater on a nominal scale. For example, variable *gender* often has two categories, while *party*
800 *affiliation* in U.S. may have four, democrat, republican, independent, and others.

801 **I.1.4. Distribution Skew (s_k).** Distribution, aka base rate, frequency, marginal, or
802 prevalence, refers to the pattern of percentage occurrences, e.g. 49% female and 51% male,
803 or 5% unhealthy and 95% healthy (27,29,76,77,53–58,65,75). Major indices are symmetrical,
804 centered on 50&50% distribution. Accordingly, this study folded the original distribution to
805 create *distribution skew* (s_k), which served as a main independent variable.

806 **I.1.5. Difficulty (df).** *Difficulty* (df) represents the combination of all factors that make
807 rating inaccurate, including 1) *task difficulty*: Some tasks are more difficult than others; 2)
808 *rater difficulty*: Some raters are less capable, focused, or motivated than others, which
809 increases difficulty; 3) *instrument difficulty*: Instruments are means that help raters to
810 accomplish a task, including organization, instruction, training, and equipment. Deficient
811 instruments increase difficulty. This study fixed instrument difficulty at the lower end by
812 giving easily understood tasks and instructions. We manipulated task difficulty and assumed
813 variation in rater difficulty.

814 **I.2. Five Viewpoints**

815 Five viewpoints have influenced experts' understanding of interrater reliability. They
816 are also the theoretical focal points of this study.

817 **I.2.1. Chance agreement inflates a_o .** In the academic literature on interrater-
818 interrater reliability, likely the earliest and the most widely received viewpoint is that percent
819 agreement (a_o) inflates reliability by overlooking *chance agreement* (a_c). Consequently a_o is

820 considered “the most primitive,” (2 p38) “inadequate,” (13 pp187&193) and “flawed,” (12
821 p80) therefore “should *not* be used.” (3–5,13 p187). Removing chance agreement is the core
822 or the stated mission of early indices, e.g., Benini’s β (11). Bennett et al’s S (15), Goodman &
823 Kruskal’s λ_r (78) and Guttman’s ρ (79). Of these, only S remains in regular use today (28,
824 29).

825 **I.2.2. Rating category inflates S , I_r , and AC_I .** Another widely shared viewpoint is
826 that S depends on category while it should not. Large number of *categories*, even if empty,
827 deflates chance estimates of S (S_{ac}), thereby inflates S (16,49,56–58,66). The criticism also
828 applies to six equivalents or special cases of S , namely C (80), G (81,82), k_n (61), $PABAK$
829 (83), RE (84), and *redefined* Pi (85).

830 Perreault & Leigh took the square root of S to produce I_r (7). Gwet^{41–44} incorporated
831 the entire S into his AC_I . So category affects I_r and AC_I in a similar way as it affects S ,
832 according to mathematical analysis and simulation (14,56–58), although some consider I_r
833 “the best” (87,88,89 p.384). As I_r regularly produces higher scores than other indices, its
834 popularity has grown fast in some fields (25).

835 Eliminating category effect was a main justification for Scott (16) to offer π , which in
836 turn inspired Cohen’s κ (2) and Krippendorff’s α (19,67). Not suffering from category effect
837 is a main reason that methodologists recommend π , κ or α over alternatives (12,14).

838 **I.2.3. Distribution skew deflates π , κ & α .** Considered the “statistics of choice” (90
839 p140), κ is by far the most often used index across disciplines, followed by π and α
840 (3,4,94,14,51,56–58,91–93).

841 A controversial viewpoint is that π , κ and α depends on distribution skew while they
842 should not. The trio, critics argue, mistakenly assumes that more skewed distributions create
843 more chance agreements. Consequently, higher or lower prevalence of a variable, e.g.,
844 disease, produces larger estimates of chance agreement, thereby deflates estimated reliability
845 (3,4,56–58,61–63,65,74,76,77,5,83,93–100,7,10,27,29,51,53,55).

846 By contrast, AC_I assumes a negative skew effect on chance agreement, while I_r
847 follows S to assume no skew effect.

848 The alleged dependence of π , κ and α on skew ignited repeated and spirited debates.
849 Experts defended κ by reaffirming its validity, extending its application, or teaching its use
850 (48,75,101–108). Rogot & Goldberg introduced A_2 , a mathematical equivalent of κ (109).
851 Byrt and colleagues introduced BAK (83), and Siegel & Castellan introduced *Revised K*
852 (110), which are two equivalents of π . Krippendorff advocated and defended α vigorously
853 (12,49,50). Zwick recommended π over κ and S (10), while Hsu & Field recommended κ
854 over π (48). Vach opined that the dependence on skew is harmless (107 p655), and
855 Krippendorff acclaimed that the dependence is desirable and by design (49,50).

856 **I.2.4. Reliability indices overlook task difficulty.** An emerging viewpoint is that
857 indices of interrater reliability should depend on *task difficulty*, but they do not. More
858 difficult tasks induce more chance rating, therefore more chance agreements
859 (9,14,86,95,29,38,53–58). Krippendorff , however, opined the opposite, that “more complex”
860 tasks lead to “very small” chance agreement (50 p488).

861 **I.2.5. Indices assume intentional and maximum random rating.** Among the most
862 fundamental hence the most forcefully debated views is that the chance-adjusted indices all
863 assume intentional and maximum random rating by conspiring raters, which include all raters
864 for all ratings, all the time (9,14,23,24,26,28,52,111). The raters, according to this
865 assumption, agree *a priori* to do the following -

- 866 1) To “rate” at the commands of randomization devices, e.g., randomly thrown coins,
867 rolled dice, or drawn marbles, virtual or actual, without looking at the subjects under rating,
868 2) To rate truthfully *only* when the randomization devices disagree with each other,
869 therefore rendering no consistent command for raters to follow.

870 Krippendorff rejected this view regarding Krippendorff’s α , and characterized the
871 discussion as “strange, almost conspiratorial uses of language.” (50).

872 Bipolar all-or-nothing assumptions were detected hidden in the indices. Percent
873 agreement assumes absolutely no random rating, while the chance-adjusted indices assume

874 intentional and maximum random rating. The latter group assume that raters draw virtual or
875 actual marbles before any “rating;” they “rate” by the order of the marbles whenever the
876 marbles agree to give a consistent order; they rate honestly only when the marbles disagree
877 with each other thereby giving no consistent order (9,14,29,38,39,56–58).

878 Different indices assume different ways that raters arrange the virtual or actual
879 marbles for the random drawing and rating. S , I_r and AC_I assume that raters arrange the
880 marbles evenly across color types that are matched with rating categories, causing the triad's
881 dependence on rating category. π , κ And α assume that raters match the distribution of marble
882 colors to the pre-determined but post-reported target distribution, causing the trio's
883 dependence on target distribution and skew. As said, Krippendorff denied that α makes such
884 assumptions (50,112–114).

885 The key questions, therefore, are about rater behavior. What behaviors are assumed?
886 What behaviors take place? Do the assumptions match the behaviors? Reliability researchers
887 rely on theoretical arguments, mathematical derivation, fictitious examples, naturalistic
888 comparisons, and Monte Carlo simulation. A systematic observation of rater behavior is
889 needed to inform the debates over rater behavior.

890 This paper reports a controlled experiment that manipulated category, skew, and
891 difficulty, and observed raters' behavioral responses. Seven indices of interrater reliability
892 were tested against the observed behavior. The findings also apply to the two equivalents of

893 a_o , six equivalents of S , two equivalents of π , and one equivalent of κ , covering 18 indices in
894 total.

895 **II. Reconstructed Experiment with Golden-Standard Task**

896 Section II details the design and the execution of the reconstructed experiment that
897 provided the main empirical evidence for this study.

898 We programmed a website that asked raters to identify the longest bar from several
899 bars (Figure 1). Two of the independent variables, category, and difficulty, were manipulated
900 by programming the website.

901 **II.1. Manipulating Category (C).** *Category (C)* was manipulated by giving raters
902 two, four, six or eight bars to choose from. Thus, C had four values, 2, 4, 6 and 8.

903 **II.2. Manipulating Difficulty (d_f).** *Task difficulty (d_f)* was manipulated by varying
904 the differences between two longest bars. The differences ranged from one pixel, the smallest
905 controllable element on a computer screen, to eight pixels, which were clear to nearly
906 everyone. The variable d_f was linearly transformed to a 0~1 scale where 1 represents the most
907 difficult.

908 The two longest bars (*long bars*) were 200 pixels long plus or minus 0~4 pixels for
909 the manipulation of difficulty. The lateral distance between long bars was fixed at 150 pixels
910 to minimize distance effect.

911 We confined the main competition between the long bars. Few raters chose the short

912 bars as they were clearly shorter, which made this experiment very close to Scott's empty-
913 category assumption and minimized the correlation between *category* and *difficulty* (16).

914 **II.3. Creating One-way Golden Standard.** A *gold standard* is a consensus criterion
915 under which judgments can be made with certainty. Reliability indices are standards to
916 evaluate instruments. Now that we are to evaluate the standards, a golden standard would be
917 helpful if available. The longest-bar task provides such a golden standard. Through
918 programming codes, we the researchers always know with certainty which bar was the
919 longest, and whether each rating decision was right or wrong, based on which chance
920 agreement and true reliability can be calculated and analyzed. We use “golden standard” as a
921 stronger term than “gold standard.” The latter term was borrowed by Rudd in 1979 from
922 economics where it referred to the value of gold as a monetary standard (37,115).

923 So that variables vary, the golden standard needs to be equipped with a one-way
924 mirror that is always crystal clear to researchers, but variably clear to participants. The
925 longest-bar task also provides this figurative or virtual mirror, as the task was designed such
926 that raters sometimes knew with near certainty, but sometimes did not, thereby they had
927 opportunities to rate randomly and agree by chance.

928 **II.4. Pairing Rater Responses.** Each rater rated 10 items per period and was given
929 summary statistics of right and wrong at the end of each period. The task was made to
930 resemble an online game or IQ test to maintain raters' attention and focus. Items per period

931 were limited to 10 to reduce clutter effect (116,117). Number of bars, level of difficulty, and
932 the location of long bars were randomly rotated to minimize the effects of learning, fatigue,
933 boredom, serial position, rater idiosyncrasies, and other possible confounders (117–121).

934 The same 10 items were rated again in the same order by the next rater available.
935 After completing 10 items, a rater may choose to rate 10 more. He or she might be given 10
936 unpaired items rated by another rater, or 10 new items if all rated items had been paired. The
937 process repeated until the end of data collection.

938 The data collection took place in a three-month period. Students, teachers,
939 researchers, technicians, managers, office workers and other professionals from 15 colleges
940 and two research firms in America, China mainland, Hong Kong, Macau and Singapore
941 participated as a part of their class exercises, professional training, or work assignments.
942 They registered 383 web names and logged on from 53 Asian, European and North American
943 cities. They rated a total of 22,290 items, of which 19,900 were successfully paired,
944 producing 9,950 paired responses, from which we sampled and resampled to reconstruct 384
945 rating sessions to form a between-subject (session) experiment that we report below.

946 **II.5. Manipulating Skew (s_k).** As the longest bar is either at the left or right side of
947 the second longest bar, we defined *distribution* as the left-and-right percentage. For example,
948 when 1% of the rated screens had the longest bar at the left, the distribution is denoted 1&99.
949 Five levels were chosen: 1&99, 25&75, 50&50, 75&25, and 99&1, the last of which

950 represented 99% left & 1% right. 0&100 and 100&0 were omitted as π , κ and α would be
951 undefined.

952 It is skew, but not the unfolded distribution, that's expected to affect the indices
953 (14,29,61,65,76). Therefore, *skew* (s_k) was operationalized as *distribution folded in the*
954 *middle*. 1&99 and 99&1 were both assigned $s_k=0.99$, for the highest skew. 50&50 was
955 assigned $s_k=0.5$ for the lowest skew, and 25&75 and 75&25 were both assigned $s_k=0.75$ for
956 moderate skew. Variable *skew* (s_k) ranged 0.5~0.99.

957 **II.6. Reconstructing Rating Sessions.** To reconstruct the first rating session, we
958 randomly sampled without replacement 100 paired rating responses ($N_t=100$) requiring two
959 *categories* ($C=2$), lowest *difficulty* ($d_f=0$), and highest *skew* ($s_k=.99$). After recording the
960 variable and response information, we returned the sample to the population of 9,950.

961 To reconstruct the second rating session, we drew another random sample of 100 pairs
962 requiring four *categories* ($C=4$) while the other two variables, *difficulty* and *skew*, remained
963 $d_f=0$ and $s_k=.99$. Again, we returned each pair back to the population after recording the
964 needed information. We then reconstructed the third session, then the fourth, and so on. We
965 repeated the process for every combination of *category*, *difficulty*, and *skew*, producing
966 $4*8*3=96$ sessions.

967 A few cases can significantly affect π , κ and α when distribution is skewed (51,56–
968 58,62,63,77,93,108,122). To assure stable effects, we resampled three more times to

969 quadruple the number of sessions, so $N_c=96*4=384$, which was the total number of the
970 reconstructed rating sessions that constituted the “subjects” for this experiment. Each skew
971 condition had an equal number of high- and low- prevalence sessions, that is, each skew=.99
972 condition had two 1&99 sessions and two 99&1 sessions, and each skew=.75 condition had
973 two 25%75 conditions and two 75&25 sessions.

974 **II.7. Reconstructed Experiment in Summary.** This was a 4X8X3 between-subject
975 controlled experiment with 4 subjects per cell where each subject was a rating session, as
976 shown in Table 1. The execution took two stages. The first was *individual-level treatment-*
977 *response*, during which individual-level independent variables, category and difficulty, were
978 manipulated, stimulus and treatment were administered, and individual responses were
979 recorded. The second was *group-level reconstruction*, during which individual responses
980 were sampled and resampled, and the group-level independent variable, skew, was
981 manipulated.

982 While the treatment and response collection followed the procedure of typical
983 controlled experiment (36), the sampling and resampling benefited from the theories and
984 techniques of bootstrap (32,33); jackknife (34) and Monte Carlo simulation (35).

985 Simulation is a powerful tool for understanding reliability. But simulations do not
986 measure behavior. They presume certain behaviors then examine their consequences (53,55–
987 58,123). A typical individual-level experiment is unsuitable because reliability indices are

988 meaningful only for rating sessions. A session-level experiment would require hundreds of
989 rating sessions, which would be too costly and too difficult to administer. Each rating session
990 would require a fixed level for each independent variable, e.g., all tasks are extremely
991 difficult, have eight categories, and 99% are left, which would deviate too much from
992 realistic rating. Reconstructed experiment offers a useful and feasible addition to our toolkit,
993 allowing observed rater behaviors to be factored into the debate over how raters behave.

994 **III. Variable Measurements and Calculations**

995 **III.1. Calculating Chance Agreement (o_{ac}).** The raters reported few agreements on
996 short bars (0.45%, Table 2), confirming that the main competition was successfully limited
997 between the long bars. It also simplifies the calculation for chance agreement. Assuming no
998 deliberate and systematic errors, each *erroneous agreement* (o_{ae}), the agreement between two
999 raters choosing a same wrong bar, is considered random. Because there were only two real
1000 choices, the probability theory predicates an equal number of agreements falling on the
1001 longest bars, thus being correct by chance. Therefore, *observed chance agreement* (o_{ac}) was
1002 calculated by doubling the directly observed erroneous agreement o_{ae} :

$$o_{ac} = 2 * o_{ae} \quad (3)$$

1003 To be sure, we derived another formula for o_{ac} assuming that sometimes raters had
1004 four, six, or eight real choices, as described in Section III.2 below. The two measures yielded
1005 essentially the same results. As Eq. 3 is simpler and easier to trace back to the directly

1006 observed o_{ae} , we report statistics based on Eq. 3.

1007 **III.2. Alternative Calculation of Observed Chance Agreement (o_{ac}).** We identified
1008 two formulas for calculating the observed chance agreement (o_{ac}). The findings section of the
1009 manuscript reports the results based on the simpler formula (Eq. 3). All analyses involving o_{ac}
1010 were performed twice using the two different formulas, which produced essentially the same
1011 results. We describe the alternative formula (Eq. 4) below.

1012 Some agreements are right, some are erroneous. This study directly observed
1013 erroneous agreement (o_{ae}). As we assume no systematic error, all o_{ae} are assumed to have
1014 come from chance rating, which constitutes the first part of the chance agreement to be
1015 estimated.

1016 The observed right agreement (o_{ar}) includes randomly and systematically right
1017 agreement. We need to estimate the former. Due to our design of two long bars and several (0,
1018 2, 4, 6) short bars, the chance agreement came from two types of random selection: between
1019 two long bars, and among all bars. When the latter results in an agreement on the longest bar,
1020 we call it *right agreement from random choices among all bars* (a_{ra})

1021 All agreement on the short bars resulted from raters choosing randomly among all
1022 bars. With C categories, $1/C$ of such random choices should fall on each bar, including the
1023 longest bar. Suppose there are four bars ($C=4$), and o_{s4} represents observed agreement on the
1024 two short bars, the right agreement (on the longest bar) from choosing randomly among four
1025 bars equals the agreement on each short bar, which is $o_{s4}/2$. Similarly, the right agreement

1026 from choosing randomly among six or eight bars is $o_{s6}/4$ or $o_{s8}/6$, respectively. So the total
1027 amount of right agreement from random selection among all bars is
1028 $a_{ra}=(o_{s4}/2)+(o_{s6}/4)+(o_{s8}/6)$, which constitutes the second part of the chance agreement we want
1029 to estimate.

1030 Of all observed agreements on the second longest bar (o_{a2}), some came from random
1031 selection among all bars (a_{ra}), and the rest ($o_{a2}-a_{ra}$) came from random selection between the
1032 two long bars. The same amount ($o_{a2}-a_{ra}$) should fall on the longest bar, which constitutes the
1033 last part of the chance agreement we want to estimate.

1034 Adding up the three parts, the observed chance agreement o_{ac} is:

$$o_{ac} = o_{ae} + a_{ra} + (o_{a2} - a_{ra}) = o_{ae} + o_{a2} \quad (4)$$

1035 As mentioned, the two approaches of calculating o_{ac} produced very small differences
1036 in means and even smaller differences in correlations. The two formulas therefore corroborate
1037 each other.

1038 **III.3. Calculating Observed Reliability (o_{ri}).** Observed reliability (o_{ri}) is observed
1039 agreement (a_o) minus observed chance agreement (o_{ac}):

$$o_{ri} = a_o - o_{ac} \quad (5)$$

1040 **IV. Statistical Indicators**

1041 Typical studies calculate estimators to estimate estimands, the targets of estimations.

1042 This study observed estimands to evaluate their estimators. We adopted and adapted common

1043 indicators, *mean*, *error*, and r^2 , to analyze data from this novel design with novel objectives.

1044 To guide our choices, we first review the two functions of interrater reliability as estimators.

1045 **IV.1. Approximating and predictive functions of reliability indices.** Reliability

1046 indices serve two functions. One is to compare an instrument with fixed benchmarks, such as

1047 0 for absence of reliability, 0.67 for highly tentative reliability, 0.8 for acceptable reliability,

1048 and 1 for perfect reliability (19 p147). This function requires an index to *approximate* true

1049 reliability in order to *place* accurate scores on instruments, and we need a proximity

1050 measure(s) to assess and analyze indices' ability to approximate true reliability.

1051 Another function is to compare instruments with each other in order to *differentiate*

1052 them. This function requires an index to accurately *predict* true reliability, which means to be

1053 highly and positively correlated with its estimation target, so that it almost always gives

1054 higher scores to more reliable instruments and lower scores to less reliable instruments. We

1055 need a correlational measure(s) to evaluate the indices' ability to predict true reliability.

1056 If an index always approximates the reliability of every individual session perfectly, it

1057 also predicts perfectly. Assuming no perfection, however, the prediction-proximity relation is

1058 more complicated. A good predictor is not necessarily a good approximator. For example, if a

1059 perfect predictor always overestimates by a constant, it's still a perfect predictor, because all

1060 instruments benefit equally. Conversely, a good approximator is not necessarily a good

1061 predictor. While a dreadful approximator gives higher score to worse instruments and lower

1062 scores to better instruments, its errors could offset each other to make it a perfect
1063 approximator on average. Therefore, both proximity and prediction measures are needed.

1064 **IV.2. Proximity Measure I -- Error of Mean (e_m).** An intuitive proximity measure is
1065 *error of mean* (e_m), defined as the difference between the grand average (*mean*) of
1066 estimations (r_i or a_c) and the grand average (*mean*) of estimation targets (o_{ri} and o_{ac}) . For any
1067 reliability index r_i and chance estimator a_c , the error of mean (e_m) calculations are shown as
1068 Eqs. 6 and 7.

$$e_m(r_i) = \text{mean}(r_i) - \text{mean}(o_{ri}) \quad -1 \leq e_m(r_i) \leq 1 \quad (6)$$

$$e_m(a_c) = \text{mean}(a_c) - \text{mean}(o_{ac}) \quad -1 \leq e_m(a_c) \leq 1 \quad (7)$$

1069 For example, the difference ($e_m(r_i)$) between κ estimation (r_i) and observed reliability
1070 (o_{ri}), averaged across 384 sessions, would indicate one aspect of κ 's inaccuracy.

1071 As a vector, a positive e_m indicates overestimation, while a negative e_m indicates
1072 underestimation. A near zero e_m , however, does not necessarily indicate accuracy for
1073 individual rating sessions. Overestimations and underestimations of individual sessions may
1074 offset each other to create a small e_m , a phenomenon known as *aggregation bias* or *ecological*
1075 *fallacy* (124,125).

1076 In typical research, however, overestimation of one study does not offset the
1077 underestimation of another study. Errors of all directions accumulate or even multiply in
1078 terms of social impact. We need an additional measure, which is described below.

1079 **IV.3. Proximity Measure II -- Mean of Errors (m_e).** To avoid aggregation bias, we
 1080 took the absolute value of the estimation error of each session, $|r_i - o_{ri}|$ and $|a_c - o_{ac}|$, and
 1081 averaged them across all 384 sessions. The results are *mean of errors* (m_e) for reliability and
 1082 chance estimations for reliability (r_i) and chance errors (a_c), as shown in Eqs. 8 & 9:

$$\mathbf{m}_e(\mathbf{r}_i) = \mathbf{mean}(|\mathbf{r}_i - \mathbf{o}_{ri}|) \quad 0 \leq m_e(r_i) \leq 1 \quad (8)$$

$$\mathbf{m}_e(\mathbf{a}_c) = \mathbf{mean}(|\mathbf{a}_c - \mathbf{o}_{ac}|) \quad 0 \leq m_e(a_c) \leq 1 \quad (9)$$

1083 Smaller m_e indicates a smaller error hence a better estimator. As a scalar, however,
 1084 m_e does not differentiate overestimations from underestimations, which vector e_m does.

1085 The spreads of our main variables varied significantly (Lines 4,5,10 &11 of Table 3),
 1086 which presents another concern. A narrower spread makes e_m and m_e look closer to zero
 1087 because their baselines (-1~1 or 0~1) do not change with spreads, producing a statistical
 1088 version of *baseline bias* (126) or *scale of reference bias* (127).

1089 **IV.4. Predictive Accuracy and Share of Influence -- Directional r^2 (dr^2).** As a ratio
 1090 of regression prediction over total variance, r^2 is commonly used to measure *predictive*
 1091 *accuracy* (128–130). As a percent of dependent variance explained by independent
 1092 variable(s), r^2 also indicates *share of influence* (128,129,131). As a scalar, however, r^2 does
 1093 not signal direction, while direction is important for this study. There are conflicting
 1094 expectations about how difficulty or skew affects chance agreement, for example. We added
 1095 the sign of r to r^2 to produce a *directional r squared* (dr^2):

$$\mathbf{dr}^2 = \mathbf{r} * |\mathbf{r}| \quad - 1 \leq \mathbf{dr}^2 \leq 1 \quad (10)$$

1096 We use dr^2 as the main indicator of indices' predictive accuracy and various
1097 variables' share of influence.

1098 **IV.5. Regression vs ANOVA.** Experimenters often employ ANOVA for analyzing
1099 data. The independent variables of this experiment are on ratio scales, which can be more
1100 efficiently analyzed with regression. As regression and ANOVA are mathematically
1101 equivalent, there is no loss in essential information or accuracy.

1102 **V. Benchmarks and Thresholds**

1103 This is the first time interrater reliability estimators and their chance agreement
1104 estimators are evaluated against their respective estimands, the observed true reliability and
1105 observed true chance agreement. No preestablished benchmarks or thresholds are available.
1106 Before reporting the outcome, this section lays out the principles that guide the evaluation.
1107 Besides helping the reviewers to evaluate our evaluation, we also hope that explicating the
1108 principles, if published, may start a conversation about what criteria and principles are
1109 appropriate for this type of evaluations.

1110 **V.1. Ideal index outperforms all others.** An ideal index outperforms all other
1111 indices on all indicators, producing the largest dr^2 and smallest m_e and e_m for both reliability
1112 and chance estimations. Since no such index emerged, the following principles applied.

1113 **V.2. Reliability over chance agreement.** While chance estimation is important for
1114 understanding an index's inside, an index's value is ultimately judged by the accuracy of its
1115 reliability estimation.

1116 **V.3. Prediction (dr^2) over approximation (m_e & e_m).** As said, a good predictor
1117 usually gives more reliable instruments higher scores, and less reliable instruments lower
1118 scores. A good predictor can be a poor approximator only when its estimations deviate from
1119 the true reliability by a near constant across all studies. If the constant can be estimated, such
1120 as in studies like this, researchers can add the constant to the estimations to improve the
1121 approximation. If the constant cannot be estimated, researchers may collectively adjust the
1122 benchmarks to reduce the impact of the across-the-board miss-approximation.

1123 When a good approximator is a poor predictor, its consequences are more severe and
1124 harder to remedy. A poor predictor often gives more reliable instruments lower scores, and
1125 less reliable instruments higher scores. A poor predictor can be a good approximator only
1126 when its errors on individual studies offset each other to lower the across-study errors. The
1127 offsetting through averaging does not remedy the underlying cause of the large estimation
1128 errors shown in the low correlation.

1129 If we cannot have both, we would trade approximating precision for differentiating
1130 precision. When evaluating reliability indices, therefore, more weights should be placed on
1131 dr^2 than m_e or e_m .

1132 **V.4. m_e over e_m .** To evaluate the indices' approximation accuracy, we place more
1133 weights on mean of errors (m_e) because it is less influenced by aggregation bias.

1134 **V.5. Primary Requirement.** Some disciplines honor $r_i > 0.8$ as the criterion for
1135 acknowledging reliability, and $r_i > 0.67$ for highly tentative acknowledgment (19,49,132).
1136 Without more reasonable precedents to following, this study tentatively adopts 0.8 and 0.67
1137 as thresholds for dr^2 , m_e and e_m . In accordance with Reasoning VI.3 above, we consider
1138 Inequality 11 a primary requirement for accepting an index's validity, where $dr^2_{(ori&ri)}$
1139 represents directional r^2 between observed reliability (o_{ri}) and an index's estimated reliability
1140 (r_i):

$$dr^2_{(ori&ri)} > 0.8 \quad - 1 \leq dr^2 \leq 1 \quad (11)$$

1141 The stated mission of chance-adjusted indices is to outperform percent agreement
1142 (a_o), which requires Inequality 12, where $dr^2_{(ori&a_o)}$ represents directional r^2 between o_{ri} and
1143 a_o .

$$dr^2_{(ori&ri)} \geq dr^2_{(ori&a_o)} \quad - 1 \leq dr^2 \leq 1 \quad (12)$$

1144 Inequality 11 applies when $dr^2_{(ori&a_o)} < 0.8$, otherwise Inequality 12 applies.

1145 **V.6. Secondary Requirement.** Inequalities 13 & 14 serve as the secondary
1146 requirement, where $m_e (r_i)$ and $m_e (a_o)$ represent respectively approximation errors (m_e) of an
1147 index (r_i) and a_o .

$$m_{e(r_i)} < 0.2 \quad 0 \leq m_e \leq 1 \quad (13)$$

$$m_{e(r_i)} \leq m_{e(a_o)} \quad 0 \leq m_e \leq 1 \quad (14)$$

1148 Inequality 13 applies when $m_{e(a_o)} > 0.2$; Inequality 14 applies otherwise. The threshold
 1149 0.2 in Ineq. 13 comes from $1 - 0.8 = 0.2$, where 0.8 is borrowed from, again, from
 1150 Krippendorff's criteria (19,49,132).

1151 **V.7. Tentative Requirement.** In case no index meets the primary and secondary
 1152 requirements, thresholds of 0.67 for dr^2 and 0.33 for m_e may be applied for tentative
 1153 acceptance, again borrowing Krippendorff's criteria (19,49,132).

1154 **V.8. Competitive requirement.** To be among the recommended, an index also needs
 1155 to outperform all other indices on at least one of the major indicators.

1156

1157

References

- 1158 1. Artstein R. Inter-annotator agreement. In: Ide N, editor. Handbook of Linguistic
1159 Annotation [Internet]. Springer Netherlands; 2017 [cited 2022 Jan 15]. p. 297–313.
1160 Available from: https://link.springer.com/chapter/10.1007/978-94-024-0881-2_11
- 1161 2. Cohen J. A coefficient of agreement for nominal scales. Educ Psychol Meas [Internet].
1162 1960 [cited 2022 Jan 15];20(1):37–46. Available from:
1163 <http://psycnet.apa.org/index.cfm?fa=search.displayRecord&uid=1960-06759-001>
- 1164 3. Feng GC. Estimating intercoder reliability: a structural equation modeling approach.
1165 Qual Quant [Internet]. 2014 Jul 20 [cited 2022 Jan 15];48(4):2355–69. Available from:
1166 <http://link.springer.com/10.1007/s11135-014-0034-7>
- 1167 4. Feng GC. Intercoder reliability indices: Disuse, misuse, and abuse. Qual Quant
1168 [Internet]. 2014 [cited 2022 Jan 15];48(3):1803–15. Available from:
1169 <http://link.springer.com/article/10.1007/s11135-013-9956-8>
- 1170 5. Feng GC. Mistakes and how to avoid mistakes in using intercoder reliability indices.
1171 Methodology [Internet]. 2015 [cited 2022 Jan 15];11(1):13–22. Available from:
1172 <http://econtent.hogrefe.com/doi/full/10.1027/1614-2241/a000086>
- 1173 6. Grayson K, Rust R. Interrater reliability. J Consum Psychol [Internet]. 2001 [cited
1174 2022 Jan 15];10(1/2):71–3. Available from:

- 1175 <http://www.ncbi.nlm.nih.gov/pubmed/22114173>
1176 <http://linkinghub.elsevier.com/retrieve/pii/S1057740801702471>
- 1177 7. Perreault WD, Leigh LE. Reliability of nominal data based on qualitative judgments. *J*
1178 *Mark Res.* 1989;26(2):135–48.
- 1179 8. Popping R. On agreement indices for nominal data. In: Saris WE, Gallhofer IN,
1180 editors. *Sociometric research: Volume I, data collection and scaling* [Internet]. 1st ed.
1181 New York, NY: St. Martin's / Springer; 1988 [cited 2022 Jan 15]. p. 90–105.
1182 Available from: http://link.springer.com/chapter/10.1007/978-1-349-19051-5_6
- 1183 9. Riffe D, Lacy S, Fico FG. *Analyzing Media Messages: Using Quantitative Content*
1184 *Analysis in Research* [Internet]. 2nd ed. Mahwah, New Jersey and London, New
1185 Jersey and London: Lawrence Erlbaum Associates, Publishers; 2005 [cited 2022 Jan
1186 15]. Available from:
1187 https://books.google.com.hk/books?hl=en&lr=&id=enCRAgAAQBAJ&oi=fnd&pg=PP1&ots=B00EbKHtj7&sig=e_EdXbsENFS9VfNJR62OrQ00_MM&redir_esc=y#v=onepage&q&f=false
- 1190 10. Zwick R. Another look at interrater agreement. *Psychol Bull* [Internet].
1191 1988;103(3):374–8. Available from: [http://www.scopus.com/inward/record.url?eid=2-](http://www.scopus.com/inward/record.url?eid=2-s2.0-0024005773&partnerID=tZOtx3y1)
1192 [s2.0-0024005773&partnerID=tZOtx3y1](http://www.scopus.com/inward/record.url?eid=2-s2.0-0024005773&partnerID=tZOtx3y1)

- 1193 11. Benini R. Principii di Demographia: Manuali Barbera Di Scienze Giuridiche Sociali e
1194 Politiche (No. 29)[Principles of demographics (Barbera Manuals of Jurisprudence and
1195 Social Policy)]. Firenze, Italy: G. Barbera; 1901.
- 1196 12. Hayes AF, Krippendorff KH. Answering the call for a standard reliability measure for
1197 coding data. *Commun Methods Meas* [Internet]. 2007 [cited 2022 Jan 15];1(1):77–89.
1198 Available from: <http://www.tandfonline.com/doi/abs/10.1080/19312450709336664>
- 1199 13. Hughes MA, Garrett DE. Intercoder reliability estimation approaches in marketing: A
1200 generalizability theory framework for quantitative data. *J Mark Res* [Internet]. 1990
1201 [cited 2022 Jan 15];27(2):185–95. Available from:
1202 [http://search.ebscohost.com/login.aspx?direct=true&db=buh&AN=9602260627&site=](http://search.ebscohost.com/login.aspx?direct=true&db=buh&AN=9602260627&site=ehost-live)
1203 [ehost-live](http://search.ebscohost.com/login.aspx?direct=true&db=buh&AN=9602260627&site=ehost-live)
- 1204 14. Zhao X, Liu JS, Deng K. Assumptions behind intercoder reliability indices. *Ann Int*
1205 *Commun Assoc* [Internet]. 2013;36(1):419–80. Available from:
1206 [http://www.tandfonline.com/doi/abs/10.1080/23808985.2013.11679142?journalCode=](http://www.tandfonline.com/doi/abs/10.1080/23808985.2013.11679142?journalCode=rica20)
1207 [rica20](http://www.tandfonline.com/doi/abs/10.1080/23808985.2013.11679142?journalCode=rica20)
- 1208 15. Bennett EM, Alpert R, Goldstein AC. Communications through limited response
1209 questioning. *Public Opin Q* [Internet]. 1954 [cited 2022 Jan 15];18:303–8. Available
1210 from:

- 1211 <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Cit>
1212 [ation&list_uids=2189948](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=2189948)
- 1213 16. Scott WA. Reliability of content analysis: The case of nominal coding. *Public Opin Q*
1214 [Internet]. 1955 [cited 2022 Jan 15];19(3):321–325. Available from:
1215 <http://www.jstor.org/stable/2746450>
- 1216 17. Krippendorff KH. Estimating the Reliability, systematic error and random error of
1217 interval data. *Educ Psychol Meas* [Internet]. 1970 [cited 2022 Jan 15];30:61–70.
1218 Available from: <http://epm.sagepub.com/content/30/1/61.short>
- 1219 18. Krippendorff KH. On generating data in communication research. *J Commun*
1220 [Internet]. 1970 [cited 2022 Jan 15];20:241–69. Available from:
1221 <http://dx.doi.org/10.1111/j.1460-2466.1970.tb00883.x>
- 1222 19. Krippendorff KH. *Content Analysis: An Introduction to its Methodology*. Thousand
1223 Oaks, CA: Sage; 1980.
- 1224 20. Button CM, Snook B, Grant MJ. Inter-rater agreement, data reliability, and the Crisis
1225 of confidence in psychological research. *Quant Methods Psychol*. 2020;16(5):467–71.
- 1226 21. Checco A, Roitero A, Maddalena E, Mizzaro S, Demartini G. Let's agree to disagree:
1227 Fixing agreement measures for crowdsourcing. *Proc Fifth AAAI Conf Hum Comput*
1228 *Crowdsourcing* [Internet]. 2017 [cited 2022 Jan 15];(Hcomp):11–20. Available from:
1229 www.aaai.org

- 1230 22. ten Hove D, Jorgensen TD, van der Ark LA. On the usefulness of interrater reliability
1231 coefficients. In: Wiberg M, Culpepper S, Janssen R, Gonzalez J, Molenaar D, editors.
1232 Quantitative Psychology: The 82nd Annual Meeting of the Psychometric Society,
1233 Zurich, Switzerland, 2017. Cham, Switzerland: Springer; 2018. p. 67–75.
- 1234 23. Zhao X. When to use Cohen’s κ , if ever? [Internet]. [Boston, USA, May,
1235 https://repository.hkbu.edu.hk/coms_conf/2/]: Paper presented at the 61st annual
1236 conference of International Communication Association.; 2011 [cited 2022 Jan 15].
1237 Available from: <https://repository.um.edu.mo/handle/10692/102423>
- 1238 24. Zhao X. When to use Scott’s π or Krippendorff’s α , if ever? [Internet]. [St. Louis,
1239 USA, August, https://repository.hkbu.edu.hk/coms_conf/3/]: Paper presented at the
1240 annual conference of Association for Education in Journalism and Mass
1241 Communication; 2011 [cited 2022 Jan 15]. Available from:
1242 <https://repository.um.edu.mo/handle/10692/102434>
- 1243 25. Zhao X, Deng K, Feng GC, Zhu L, Chan VKC. Liberal-conservative hierarchies for
1244 indices of inter-coder reliability [Internet]. Paper presented at the 62nd annual
1245 conference of International Communication Association, Phoenix, Arizona, USA,
1246 May; 2012 [cited 2022 Jan 15]. Available from:
1247 <https://repository.um.edu.mo/handle/10692/102423>

- 1248 26. Zhao X, Feng GC, Liu JS, Deng K. We agreed to measure agreement - Redefining
1249 reliability de-justifies Krippendorff's alpha. *China Media Res* [Internet]. 2018 [cited
1250 2022 Jan 15];14(2):1. Available from:
1251 <https://repository.um.edu.mo/handle/10692/25978>
- 1252 27. Conger AJ. Kappa and Rater Accuracy: Paradigms and Parameters. *Educ Psychol*
1253 *Meas* [Internet]. 2016 [cited 2022 Jan 15];0013164416663277. Available from:
1254 <http://epm.sagepub.com/content/early/2016/08/18/0013164416663277.abstract%255Cnhttp://epm.sagepub.com/content/early/2016/08/18/0013164416663277%255Cnhttp://epm.sagepub.com/content/early/2016/08/18/0013164416663277.full.pdf>
- 1255
1256
- 1257 28. Delgado R, Tibau XA. Why Cohen's Kappa should be avoided as performance
1258 measure in classification. *PLoS One* [Internet]. 2019 [cited 2022 Jan 15];14(9):1–26.
1259 Available from: <http://dx.doi.org/10.1371/journal.pone.0222916>
- 1260 29. Grove WM, Andreasen NC, McDonald-Scott P, Keller MB, Shapiro RW. Reliability
1261 studies of psychiatric diagnosis: Theory and practice. *Arch Gen Psychiatry*.
1262 1981;38(4):408–13.
- 1263 30. Riffe D, Lacy S, Fico FG, Watson B. Analyzing media messages: Using quantitative
1264 content analysis in research (4th ed.). [Internet]. 4th ed. New York: Routledge; 2019
1265 [cited 2022 Jan 15]. Available from:
1266 <https://www.taylorfrancis.com/books/9780429464287>

- 1267 31. Zhao X. Selective spiral ——A mega, meta, predictive and presumptive theory of
1268 communication. Chinese J Journal Commun [Internet]. 2018 [cited 2022 Jan
1269 15];40(2):140–53. Available from: <http://cjjc.ruc.edu.cn/EN/Y2018/V40/I2/140>
- 1270 32. Efron B. Bootstrap Methods: Another Look at the Jackknife. Ann Stat [Internet]. 1979
1271 [cited 2022 Jan 15];7(1):1–26. Available from:
1272 <http://projecteuclid.org/euclid.aos/1176344552>
1273 <https://projecteuclid.org/euclid.aos/1176344552>
- 1274 33. Efron B, Robert J. Tibshirani. An Introduction to the Bootstrap [Internet]. New York
1275 and London: Chapman & Hall; 1993 [cited 2022 Jan 15]. 257 p. Available from:
1276 <http://books.google.com/books?id=gLlpIUxRntoC&pgis=1>
- 1277 34. Shao J, Tu D. The Jackknife and Bootstrap [Internet]. Springer Series in Statistics.
1278 New York: Springer Science & Business Media; 1995 [cited 2022 Jan 15]. 516 p.
1279 Available from: <http://www.loc.gov/catdir/enhancements/fy0815/95015074-d.html>
- 1280 35. Liu JS. Monte Carlo strategies in scientific computing. New York: Springer; 2001.
- 1281 36. Montgomery DC. Design and Analysis of Experiments, 7th Edition. John Wiley &
1282 Sons. 2009.
- 1283 37. Claassen JAHR. The gold standard: not a golden standard. BMJ [Internet]. 2005 [cited
1284 2022 Jan 15];330(7500):1121. Available from: bmj.com

- 1285 38. Riffe D, Lacy S, Fico FG. Analyzing Media Messages: Using Quantitative Content
1286 Analysis in Research. Mahwah, N J: Lawrence Erlbaum Associates; 1998.
- 1287 39. Riffe D, Lacy S, Fico FG. Analyzing Media Messages : Using Quantitative Content
1288 Analysis in Research. 3rd ed. New York: Routledge; 2014.
- 1289 40. Wasserstein RL, Lazar NA. The ASA’s statement on p -Values: context, process, and
1290 purpose. Am Stat [Internet]. 2016 [cited 2022 Jan 15];70(2):129–33. Available from:
1291 <https://www.tandfonline.com/doi/full/10.1080/00031305.2016.1154108>
- 1292 41. Amrhein V, Greenland S, McShane B, others. Retire statistical significance. Nature.
1293 2019 Mar 21;567:305–7.
- 1294 42. Wasserstein RL, Schirm AL, Lazar NA. Moving to a world beyond “ $p < 0.05$.” Am
1295 Stat. 2019;73(sup1):1–19.
- 1296 43. Wilkinson L, APA Task Force. Statistical methods in psychology journals: Guidelines
1297 and explanations (Report by Task Force on Statistical Inference, APA Board of
1298 Scientific Affairs). Am Psychol [Internet]. 1999 [cited 2022 Jan 15];54(8):594–604.
1299 Available from: <http://psycnet.apa.org/journals/amp/54/8/594/>
- 1300 44. Lazar N. Time to say goodbye to “statistically significant” and embrace uncertainty,
1301 say statisticians. Retraction Watch [Internet]. 2019 Mar 21 [cited 2022 Jan 15];
1302 Available from: [https://retractionwatch.com/2019/03/21/time-to-say-goodbye-to-
1303 statistically-significant-and-embrace-uncertainty-say-statisticians/](https://retractionwatch.com/2019/03/21/time-to-say-goodbye-to-statistically-significant-and-embrace-uncertainty-say-statisticians/)

- 1304 45. Liu PL, Zhao X, Wan B. COVID-19 information exposure and vaccine hesitancy: The
1305 influence of trust in government and vaccine confidence. *Psychol Heal Med* [Internet].
1306 2021 [cited 2022 Jan 15];00(00):1–10. Available from:
1307 <https://doi.org/10.1080/13548506.2021.2014910>
- 1308 46. Zhao X. Four functions of statistical significance tests [Internet]. Presentation at the
1309 School of Statistics and Center for Data Sciences Beijing Normal University, 25th
1310 December.; 2016 [cited 2022 Jan 15]. Available from:
1311 <https://repository.um.edu.mo/handle/10692/95184>
- 1312 47. Zhao X, Ye J, Sun S, Zhen Y, Zhang Z, Xiao Q, et al. Best Title Lengths of Online
1313 Postings for Highest Read and Relay. *Journal Commun Rev* [Internet]. 2022 [cited
1314 2022 Jul 21];75(3):5–20. Available from:
1315 <https://repository.um.edu.mo/handle/10692/95320>
- 1316 48. Hsu LM, Field R. Interrater agreement measures: Comments on Kappan, Cohen's
1317 Kappa, Scott's π , and Aickin's α . *Underst Stat*. 2003;2(3):205–19.
- 1318 49. Krippendorff KH. Reliability in content analysis: Some common misconceptions and
1319 recommendations. *Hum Commun Res*. 2004;30(3):411–33.
- 1320 50. Krippendorff KH. A dissenting view on so-called paradoxes of reliability coefficients.
1321 *Ann Int Commun Assoc* [Internet]. 2013 [cited 2022 Jan 15];36(1):481–99. Available
1322 from: <http://www.tandfonline.com/doi/pdf/10.1080/23808985.2013.11679143>

- 1323 51. Lombard M, Snyder-Duch J, Bracken CC. Content analysis in mass communication:
1324 Assessment and reporting of intercoder reliability. *Hum Commun Res* [Internet]. 2002
1325 [cited 2022 Jan 15];28(4):587–604. Available from:
1326 <http://onlinelibrary.wiley.com/doi/10.1111/j.1468-2958.2002.tb00826.x/abstract>
- 1327 52. Feng GC, Zhao X. Do not force agreement – A response to Krippendorff. *Methodol*
1328 *Eur J Res Methods Behav Soc Sci* [Internet]. 2016 [cited 2022 Jan 15];12(4):145–8.
1329 Available from: <https://repository.um.edu.mo/handle/10692/26008>
- 1330 53. Gwet KL. Computing inter-rater reliability and its variance in the presence of high
1331 agreement. *Br J Math Stat Psychol* [Internet]. 2008 [cited 2022 Jan 15];61(1):29–48.
1332 Available from: <http://doi.wiley.com/10.1348/000711006X126600>
- 1333 54. Gwet KL. *Handbook of Inter-Rater Reliability: The Definitive Guide to Measuring the*
1334 *Extent of Agreement Among Raters*. 2nd ed. Gaithersburg, MD: STATAXIS
1335 Publishing Company; 2010. 197 p.
- 1336 55. Gwet KL. Variance estimation of nominal-scale inter-rater reliability with random
1337 selection of raters. *Psychometrika* [Internet]. 2008 [cited 2016 Feb 7];73(3):407–30.
1338 Available from: <http://link.springer.com/article/10.1007/s11336-007-9054-8>
- 1339 56. Feng GC. *Indexing versus Modeling Intercoder Reliability*. Hong Kong Baptist
1340 University; 2013.

- 1341 57. Feng GC. Factors affecting intercoder reliability: A Monte Carlo experiment. *Qual*
1342 *Quant* [Internet]. 2013 [cited 2022 Jan 15];47(5):2959–82. Available from:
1343 <http://link.springer.com/article/10.1007/s11135-012-9745-9>
- 1344 58. Feng GC. Underlying determinants driving agreement among coders. *Qual Quant*.
1345 2013;47(5):2983–97.
- 1346 59. Attneave F. Psychological probability as a function of experienced frequency. *J Exp*
1347 *Psychol*. 1953;46(2):81–6.
- 1348 60. Fischhoff B, Slovic P, Lichtenstein S. Knowing with certainty: The appropriateness of
1349 extreme confidence. *J Exp Psychol Hum Percept Perform* [Internet]. 1977 [cited 2022
1350 Jan 15];3(4):552–64. Available from: <http://content.apa.org/journals/xhp/3/4/552>
- 1351 61. Brennan RL, Prediger DJ. Coefficient kappa: Some uses, misuses, and alternatives.
1352 *Educ Psychol Meas* [Internet]. 1981 [cited 2022 Jan 15];41(3):687–99. Available from:
1353 <http://journals.sagepub.com/doi/10.1177/001316448104100307>
- 1354 62. Feinstein AR, Cicchetti D V. High agreement but low Kappa: II. Resolving the
1355 paradoxes. *J Clin Epidemiol*. 1990;43(6):551–8.
- 1356 63. Feinstein AR, Cicchetti D V. High agreement but low Kappa: I. the problems of two
1357 paradoxes. *J Clin Epidemiol*. 1990;43(6):543–9.
- 1358 64. Lantz CA, Nebenzahl E. Behavior and interpretation of the Kappa statistic: Resolution
1359 of the two paradoxes. *J Clin Epidemiol*. 1996;49(4):431–4.

- 1360 65. Spitznagel EL, Helzer JE, John E. Helzer., Helzer JE. A proposed solution to the base
1361 rate problem in the kappa statistic. *Arch Gen Psychiatry* [Internet]. 1985 [cited 2022
1362 Jan 15];42(7):725–8. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/4015315>
- 1363 66. Cousineau D, Laurencelle L. An unbiased estimate of global interrater agreement.
1364 *Educ Psychol Meas* [Internet]. 2016 [cited 2022 Jan 15];0013164416654740.
1365 Available from: <http://journals.sagepub.com/doi/abs/10.1177/0013164416654740>
- 1366 67. Krippendorff KH. Estimating the reliability, systematic error and random error of
1367 interval data. *Educ Psychol Meas* [Internet]. 1970 [cited 2022 Jan 15];30(1):61–70.
1368 Available from: <http://epm.sagepub.com/cgi/doi/10.1177/001316447003000105>
- 1369 68. Cousineau D, Laurencelle L. A ratio test of interrater agreement with high specificity.
1370 *Educ Psychol Meas* [Internet]. 2015 [cited 2022 Jan 15];75(6):979–1001. Available
1371 from:
1372 <http://epm.sagepub.com/content/75/6/979.abstract?&location1=all&location2=all&row>
1373 [_operator2=and&term1a=simulation&term_operator1=and&term_operator2=and&ct](http://epm.sagepub.com/content/75/6/979.abstract?&location1=all&location2=all&row)
- 1374 69. Kirilenko AP, Stepchenkova S. Inter-Coder Agreement in One-to-Many Classification:
1375 Fuzzy Kappa. *PLoS One* [Internet]. 2016 [cited 2022 Jan 15];11(3):e0149787.
1376 Available from:
1377 <http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0149787>

- 1378 70. Meehl PE. Clinical versus statistical prediction: A theoretical analysis and a review of
1379 the evidence. Minneapolis, MN: University of Minnesota Press; 1954.
- 1380 71. Dawes RM, Faust D, Meehl. PE. Clinical Versus Actuarial Judgment. *Science* (80-).
1381 1989;243(4899):1668-1674.
- 1382 72. Chaffee SH. *Communication Concept: Explication*. Newbury Park, CA: Sage
1383 Publications, Inc.; 1991. 96 p.
- 1384 73. van Swieten JC, Koudstaal PJ, Visser MC, Schouten HJ, van Gijn J. Interobserver
1385 agreement for the assessment of handicap in stroke patients. *Stroke*. 1988;19(5):604–7.
- 1386 74. Uebersax JS. The Myth of Chance Corrected Agreement [Internet]. 2009 [cited 2012
1387 Oct 18]. Available from: <http://www.john-uebersax.com/stat/kappa2.htm>
- 1388 75. Kraemer HC. Ramifications of a population model for κ as a coefficient of reliability.
1389 *Psychometrika* [Internet]. 1979 [cited 2022 Jan 15];44(4):461–72. Available from:
1390 [http://www.scopus.com/scopus/inward/record.url?eid=2-s2.0-](http://www.scopus.com/scopus/inward/record.url?eid=2-s2.0-0001586045&partnerID=40)
1391 [0001586045&partnerID=40](http://www.scopus.com/scopus/inward/record.url?eid=2-s2.0-0001586045&partnerID=40)
- 1392 76. Shrout PE, Spitzer RL, Fleiss JL. Quantification of agreement in psychiatric diagnosis
1393 revisited. *Arch Gen Psychiatry*. 1987;44(2):172–7.
- 1394 77. von Eye A, von Eye M. On the marginal dependency of Cohen’s κ . *Eur Psychol*.
1395 2008;13(4):305–15.

- 1396 78. Goodman LA, Kruskal WH. Measures of association for cross classifications. *J Am*
1397 *Stat Assoc.* 1954 Dec;49(268):732–64.
- 1398 79. Guttman L. The test-retest reliability of qualitative data. *Psychometrika.*
1399 1946;11(2):81–95.
- 1400 80. Janson S, Vegelius J. On generalizations of the G index and the Phi coefficient to
1401 nominal scales. *Multivariate Behav Res.* 1979;14(2):255–69.
- 1402 81. Guilford JP. Preparation of item scores for correlation between individuals in a Q
1403 factor analysis. Paper Presented at the Annual Convention of the Society of
1404 Multivariate Experimental Psychologists; 1961.
- 1405 82. Holley JW, Guilford JP. A note on the G-index of agreement. *Educ Psychol Meas.*
1406 1964;24(4):749–53.
- 1407 83. Byrt T, Bishop J, Carlin JB. Bias, prevalence and kappa. *J Clin Epidemiol [Internet].*
1408 1993 [cited 2022 Jan 15];46(5):423–9. Available from:
1409 <http://www.sciencedirect.com/science/article/pii/089543569390018V>
- 1410 84. Maxwell AE. Coefficients of agreement between observers and their interpretation. *Br*
1411 *J Psychiatry.* 1977;130(1):79–83.
- 1412 85. Potter WJ, Levine-Donnerstein D. Rethinking validity and reliability in content
1413 analysis. *J Appl Commun Res [Internet].* 1999 [cited 2022 Jan 15];27(3):258–84.
1414 Available from: <http://www.tandfonline.com/doi/abs/10.1080/00909889909365539>

- 1415 86. Gwet KL. Handbook of Inter-Rater Reliability: The Definitive Guide to Measuring the
1416 Extent of Agreement Among Multiple Raters [Internet]. 3rd ed. Gaithersburg, MD,
1417 USA: Advanced Analytics, LLC; 2012 [cited 2022 Jan 15]. 197 p. Available from:
1418 [https://books.google.com/books?hl=en&lr=&id=fac9BQAAQBAJ&oi=fnd&pg=PP1&](https://books.google.com/books?hl=en&lr=&id=fac9BQAAQBAJ&oi=fnd&pg=PP1&dq=Gwet+K+L&ots=UUdriDAp0a&sig=mKjbb_IW1eNG474Cb0Omp3n5BMk)
1419 [dq=Gwet+K+L&ots=UUdriDAp0a&sig=mKjbb_IW1eNG474Cb0Omp3n5BMk](https://books.google.com/books?hl=en&lr=&id=fac9BQAAQBAJ&oi=fnd&pg=PP1&dq=Gwet+K+L&ots=UUdriDAp0a&sig=mKjbb_IW1eNG474Cb0Omp3n5BMk)
- 1420 87. Ji MF, McNeal JU. How chinese children's commercials differ from those of the
1421 united states: A content analysis. J Advert [Internet]. 2001 [cited 2022 Jan
1422 15];30(3):79–92. Available from:
1423 [http://web.ebscohost.com/ehost/detail?hid=106&sid=c0d4783a-f726-4eea-9dc2-](http://web.ebscohost.com/ehost/detail?hid=106&sid=c0d4783a-f726-4eea-9dc2-4022b157f163@sessionmgr112&vid=3&bdata=JnNpdGU9ZWwhvc3QtbGl2ZQ==#db=buh&AN=5507388)
1424 [4022b157f163@sessionmgr112&vid=3&bdata=JnNpdGU9ZWwhvc3QtbGl2ZQ==#db](http://web.ebscohost.com/ehost/detail?hid=106&sid=c0d4783a-f726-4eea-9dc2-4022b157f163@sessionmgr112&vid=3&bdata=JnNpdGU9ZWwhvc3QtbGl2ZQ==#db=buh&AN=5507388)
1425 [=buh&AN=5507388](http://web.ebscohost.com/ehost/detail?hid=106&sid=c0d4783a-f726-4eea-9dc2-4022b157f163@sessionmgr112&vid=3&bdata=JnNpdGU9ZWwhvc3QtbGl2ZQ==#db=buh&AN=5507388)
- 1426 88. Kolbe RH, Burnett MS. Content-analysis research: An examination of applications
1427 with directives for improving research reliability and objectivity. J Consum Res.
1428 1991;18(2):243–50.
- 1429 89. Okazaki S, Rivas JA. A content analysis of multinationals' Web communication
1430 strategies: cross-cultural research framework and pre-testing. Internet Res.
1431 2002;12(5):380–90.
- 1432 90. Uebersax JS. Diversity of decision-making models and the measurement of interrater
1433 agreement. Psychol Bull. 1987;101(1):140–6.

- 1434 91. Andsager JL, Schwartz J. Explicating time: toward making content analysis research
1435 describing time frames more meaningful. [Chicago]: Paper Presented at Annual
1436 Conference of Association for Education in Journalism and Mass Communication,
1437 Chicago; 2012.
- 1438 92. Bakeman R. Behavioral observation and coding. *Handb Res methods Soc Personal*
1439 *Psychol* [Internet]. 2000 [cited 2022 Jan 15];138–59. Available from:
1440 [http://search.ebscohost.com/login.aspx?direct=true&db=psyh&AN=2000-07611-](http://search.ebscohost.com/login.aspx?direct=true&db=psyh&AN=2000-07611-006&site=ehost-live)
1441 [006&site=ehost-live](http://search.ebscohost.com/login.aspx?direct=true&db=psyh&AN=2000-07611-006&site=ehost-live)
- 1442 93. von Eye A. An alternative to Cohen’s kappa. Vol. 11, *European Psychologist*. 2006. p.
1443 12–24.
- 1444 94. Warrens MJ. On marginal dependencies of the 2×2 kappa. *Adv Stat* [Internet]. 2014
1445 [cited 2022 Jan 15];2014:1–6. Available from:
1446 <http://www.hindawi.com/archive/2014/759527/>
- 1447 95. Aickin M. Maximum likelihood estimation of agreement in the constant predictive
1448 probability model, and its relation to Cohen’s kappa. *Biometrics* [Internet]. 1990 [cited
1449 2022 Jan 15];46:293–302. Available from: <http://www.jstor.org/stable/2531434>
- 1450 96. Gwet KL. Inter-Rater Reliability: Dependency on Trait Prevalence and Marginal
1451 Homogeneity [Internet]. Gaithersburg, MD, USA; 2002 [cited 2022 Jan 15]. Available
1452 from: <http://hbanaszak.mjr.uw.edu.pl/TempTxt/smirra2.pdf>

- 1453 97. Gwet KL. Handbook of Inter-Rater Reliability: The Definitive Guide to Measuring the
1454 Extent of Agreement Among Raters [Internet]. 4th ed. Gaithersburg, MD: Advanced
1455 Analytics, LLC; 2014 [cited 2022 Jan 15]. 429 p. Available from:
1456 [https://books.google.com/books?hl=en&lr=&id=fac9BQAAQBAJ&oi=fnd&pg=PP1&](https://books.google.com/books?hl=en&lr=&id=fac9BQAAQBAJ&oi=fnd&pg=PP1&dq=Gwet+K+L&ots=UUdriDAp0a&sig=mKjbb_IW1eNG474Cb0Omp3n5BMk)
1457 [dq=Gwet+K+L&ots=UUdriDAp0a&sig=mKjbb_IW1eNG474Cb0Omp3n5BMk](https://books.google.com/books?hl=en&lr=&id=fac9BQAAQBAJ&oi=fnd&pg=PP1&dq=Gwet+K+L&ots=UUdriDAp0a&sig=mKjbb_IW1eNG474Cb0Omp3n5BMk)
- 1458 98. Wongpakaran N, Wongpakaran T, Wedding D, Gwet KL. A comparison of Cohen's
1459 Kappa and Gwet's AC1 when calculating inter-rater reliability coefficients: a study
1460 conducted with personality disorder samples. BMC Med Res Methodol [Internet].
1461 2013 [cited 2022 Jan 15];13:61. Available from:
1462 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3643869&tool=pmcentrez](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3643869&tool=pmcentrez&rendertype=abstract)
1463 [&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3643869&tool=pmcentrez&rendertype=abstract)
- 1464 99. Powers DWM. The problem with Kappa. In: Proceedings of the 13th Conference of
1465 the European Chapter of the Association for Computational Linguistics. Association
1466 for Computational Linguistics; 2012. p. 345–55.
- 1467 100. Jeni LA, Cohn JF, De La Torre F. Facing imbalanced data - Recommendations for the
1468 use of performance metrics. Proc - 2013 Hum Assoc Conf Affect Comput Intell
1469 Interact ACII 2013. 2013;245–51.
- 1470 101. Bloch DA, Kraemer HC. 2 x 2 Kappa coefficients: Measures of agreement or
1471 association. Biometrics. 1989;45(1):269–87.

- 1472 102. Dewey ME. Coefficients of agreement. *Br J Psychiatry*. 1983;143(5):487–9.
- 1473 103. Feuerman M, Miller AR. Relationships between statistical measures of agreement:
1474 sensitivity, specificity and kappa. *J Eval Clin Pract* [Internet]. 2008 [cited 2022 Jan
1475 15];14(5):930–3. Available from: [http://onlinelibrary.wiley.com/doi/10.1111/j.1365-](http://onlinelibrary.wiley.com/doi/10.1111/j.1365-2753.2008.00984.x/full)
1476 [2753.2008.00984.x/full](http://onlinelibrary.wiley.com/doi/10.1111/j.1365-2753.2008.00984.x/full)
- 1477 104. Kraemer HC, Bloch DA. Kappa coefficients in epidemiology: An appraisal of a
1478 reappraisal. *J Clin Epidemiol*. 1988;41(10):959–68.
- 1479 105. Kraemer HC, Periyakoil VS, Noda A. Tutorial in Biostatistics: Kappa coefficients in
1480 medical research. *Stat Med* [Internet]. 2002 [cited 2022 Jan 15];21(14):2109–29.
1481 Available from: <http://www.ncbi.nlm.nih.gov/pubmed/12111890>
- 1482 106. Roberts C. Modelling patterns of agreement for nominal scales. *Stat Med*.
1483 2008;27(6):810–30.
- 1484 107. Williamson JM, Lipsitz SR, Amita K, Manatunga. Modeling kappa for measuring
1485 dependent categorical agreement data. *Biostatistics* [Internet]. 2000 [cited 2022 Jan
1486 15];1(2):191–202. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/12933519>
- 1487 108. Vach W. The dependence of Cohen’s kappa on the prevalence does not matter. Vol.
1488 58, *Journal of Clinical Epidemiology*. 2005. p. 655–61.
- 1489 109. Rogot E, Irving D, Goldberg. A proposed index for measuring agreement in test-retest
1490 studies. *J Chronic Dis*. 1966;19(9):991–1006.

- 1491 110. Siegel S, Castellan JNJ. Nonparametric statistics for the behavioural sciences
1492 [Internet]. 2nd ed. MacGraw Hill; 1988 [cited 2022 Jan 15]. 213–214 p. Available
1493 from:
1494 [http://scholar.google.com/scholar?hl=en&btnG=Search&q=intitle:Non+parametric+sta](http://scholar.google.com/scholar?hl=en&btnG=Search&q=intitle:Non+parametric+statistics+for+the+behavioural+sciences#9)
1495 [tistics+for+the+behavioural+sciences#9](http://scholar.google.com/scholar?hl=en&btnG=Search&q=intitle:Non+parametric+statistics+for+the+behavioural+sciences#9)
- 1496 111. Grove WM, Andreasen NC, McDonald-Scott P, Keller MB, Shapiro RW. Reliability
1497 studies of psychiatric diagnosis: Theory and practice. Arch Gen Psychiatry.
1498 1981;38(4):408–13.
- 1499 112. Krippendorff KH. Misunderstanding reliability. Methodology [Internet]. 2016 [cited
1500 2022 Jan 15];12(4):139–44. Available from:
1501 <http://econtent.hogrefe.com/doi/full/10.1027/1614-2241/a000119>
- 1502 113. Krippendorff KH. The changing landscape of content analysis: Reflections on social
1503 construction of reality and beyond. So CYK, editor. Commun Soc [Internet]. 2019
1504 [cited 2022 Jan 15];47(47):1–27. Available from:
1505 https://repository.upenn.edu/asc_papers/604
- 1506 114. Krippendorff KH. Reliability in content analysis. Hum Commun Res [Internet]. 2004
1507 [cited 2022 Jan 15];30(3):411–33. Available from:
1508 <http://onlinelibrary.wiley.com/doi/10.1111/j.1468-2958.2004.tb00738.x/abstract>

- 1509 115. Rudd P. In Search of the Gold Standard for Compliance Measurement. *Arch Intern*
1510 *Med* [Internet]. 1979 Jun 1 [cited 2022 Jan 15];139(6):627–8. Available from:
1511 <https://jamanetwork.com/journals/jamainternalmedicine/fullarticle/589465>
- 1512 116. Webb PH, Michael L. Ray. Effects of TV Clutter. *J Advert Res*. 1979;19(3):7–12.
- 1513 117. Zhao X. Clutter and serial order redefined and retested. *J Advert Res* [Internet]. 1997
1514 [cited 2022 Jan 15];37(5):57–73. Available from: [https://works.bepress.com/xinshu-](https://works.bepress.com/xinshu-zhao/11/)
1515 [zhao/11/](https://works.bepress.com/xinshu-zhao/11/)
- 1516 118. Jeong Y, Tran H, Zhao X. How much is too much? *J Advert Res*. 2012;52(1):87–101.
- 1517 119. Li C. Primacy effect or recency effect? A long-term memory test of super bowl
1518 commercials. *J Consum Behav*. 2010;9(1):32–44.
- 1519 120. Pieters RGM, Bijmolt THA. Consumer Memory for Television Advertising: A Field
1520 Study of Duration, Serial Position, and Competition Effects. *J Consum Res*.
1521 1997;23(4):362.
- 1522 121. Terry WS. Serial position effects in recall of television commercials. *J Gen Psychol*
1523 [Internet]. 2005 [cited 2022 Jan 15];132(2):151–63. Available from:
1524 <http://www.ncbi.nlm.nih.gov/pubmed/15871298>
- 1525 122. Hoehler FK. Bias and prevalence effects on kappa viewed in terms of sensitivity and
1526 specificity. *J Clin Epidemiol*. 2000;53(5):499–503.

- 1527 123. Cicchetti D V., Shoinralter D, Peter J. Tyrer. The effect of number of rating scale
1528 categories on levels of interrater reliability: A Monte Carlo investigation. *Appl Psychol*
1529 *Meas.* 1985;9(1):31–6.
- 1530 124. Achen CH, Shively WP. Cross-Level Inference. In: *Cross-Level Inference*. 1995. p. 1–
1531 29.
- 1532 125. Heerink N, Mulatu A, Bulte E, Mulatu A. Income inequality and the environment:
1533 Aggregation bias in environmental Kuznets curves. *Ecol Econ.* 2001;38(3):359–67.
- 1534 126. Kinsman RA, Staudenmayer H. Baseline levels in muscle relaxation training. *Appl*
1535 *Psychophysiol Biofeedback.* 1978;3(1):97–104.
- 1536 127. Groot W. Adaptation and scale of reference bias in self-assessments of quality of life. *J*
1537 *Health Econ.* 2000;19(3):403–20.
- 1538 128. Cohen J. *Statistical Power Analysis for the Behavioral Sciences* [Internet]. 2nd ed.
1539 Vol. 2nd. Hillsdale, New Jersey: Erihaum; 1988 [cited 2022 Jan 15]. 567 p. Available
1540 from: <http://books.google.com/books?id=TI0N2lRAO9oC&pgis=1>
- 1541 129. Cohen J. A power primer. *Psychol Bull* [Internet]. 1992 [cited 2022 Jan
1542 15];112(1):155–9. Available from: [http://doi.apa.org/getdoi.cfm?doi=10.1037/0033-](http://doi.apa.org/getdoi.cfm?doi=10.1037/0033-2909.112.1.155)
1543 [2909.112.1.155](http://doi.apa.org/getdoi.cfm?doi=10.1037/0033-2909.112.1.155)
- 1544 130. Stigler SM. Francis Galton’s Account of the Invention of Correlation. *Stat Sci.*
1545 1989;4(2):73–9.

- 1546 131. Steel RG, James H. Torrie. Principles and Procedures of Statistics. New York:
1547 McGraw-Hill; 1960.
- 1548 132. Krippendorff KH. Content Analysis: An Introduction to its Methodology. 3rd ed.
1549 Thousand Oaks, CA: Sage Publications; 2012.
- 1550 133. Gwet KL. Computing inter-rater reliability and its variance in the presence of high
1551 agreement. Br J Math Stat Psychol [Internet]. 2008 [cited 2022 Jan 15];61(1):29–48.
1552 Available from: <http://onlinelibrary.wiley.com/doi/10.1348/000711006X126600/full>
- 1553 134. Krippendorff KH. Bivariate agreement coefficients for reliability of data. Sociol
1554 Methodol [Internet]. 1970 [cited 2022 Jan 15];2:139–50. Available from:
1555 <http://www.jstor.org/stable/270787>
- 1556

Tables and Figures

Table 1 A Category (C) by Difficulty (d_f) by Skew (s_k)
- Reconstructed Experiment *

Across: <i>Distribution & Skew (s_k)</i>		50&50 $s_k=0.5$				25&75, 75&25 $s_k=0.75$				1&99, 99&1 $s_k=0.99$			
Across: <i>Category (C)</i>		2	4	6	8	2	4	6	8	2	4	6	8
difference in pixels (p_x)	<i>Difficulty</i> $d_f=(8-p_x)/7$												
1	=1.000	4	4	4	4	4	4	4	4	4	4	4	4
2	≈0.8571	4	4	4	4	4	4	4	4	4	4	4	4
3	≈0.7143	4	4	4	4	4	4	4	4	4	4	4	4
4	≈0.5714	4	4	4	4	4	4	4	4	4	4	4	4
5	≈0.4286	4	4	4	4	4	4	4	4	4	4	4	4
6	≈0.2857	4	4	4	4	4	4	4	4	4	4	4	4
7	≈0.1429	4	4	4	4	4	4	4	4	4	4	4	4
8	=0.0000	4	4	4	4	4	4	4	4	4	4	4	4

* Main cell entries are number of reconstructed rating sessions (subjects) in each experimental condition (cell).

Table 2 Concepts and Variables

		Down: Author or Origin		Reliability (True Agreement)		Chance Agreement	
		generic for any index		r_i		a_c	
Dependent Variables	Index Estimation	%Agreement (unknown author)		a_o		aO_{ac}	
		Bennett et al (1954)(15)		S		S_{ac}	
		Perreault & Leigh (1989) (7)		I_r		I_{rac}	
		Gwet (2002, 2008, 2010, 2012)(54)·(96,133)·(86)		AC_I		AC_{ac}	
		Scott (1955) (16)		π		π_{ac}	
		Cohen (1960) (2)		κ		κ_{ac}	
		Krippendorff (1970, 1980)(19,67,134)		α		α_{ac}	
	Empirical Observation	Primary Indicator		O_{ri} observed interrater reliability		O_{ac} observed chance agreement	
		Secondary Indicator (used in calculation)		O_{ar} observed right agreement		O_{ae} observed erroneous agreement	
				a_o observed agreement		d_o observed disagreement	
Independent Variables	Denotation	C		s_k		d_f or e_s	
	Concept	Category		Distribution Skew		Difficulty or Easiness	
Other Concepts	Denotation	e_m	m_e	S_{dm}	dr^2	N_c	N_d
	Concept	error of means (mean estimation minus mean target)	mean of errors (mean of differences between estimation and target)	standard deviation of an observed target of estimation (O_{ae} O_{ri})	directional r^2 ($dr^2 = r^* r $)	No. of rating sessions	No. of rating decisions within a session

Interrater Reliability Estimators Tested

Table 3 Effects of Estimation Targets, Category, Skew & Difficulty on Observed or Estimated Chance Agreement and Reliability (dr^2)

			A.	B.	C.	D.	E.	F.	G.	H.
	1	Right: Source or Author	Observation	%-agreement	<i>Bennett et al.</i>	<i>Perreault & Leigh</i>	<i>Gwet</i>	<i>Scott</i>	<i>Cohen</i>	<i>Krippendorff</i>
Effects on Intrdr Reliability Obsv & Ests	2	Right: Obsd / Estd Interrater Reliability as Dependent Variables Down: Independent Variables	o_{ri}	a_o	S	I_r	AC_1	π	κ	α
	3	Observed Reliability (o_{ri})	1.00***	.841***	.691***	.599***	.721***	.312***	.312***	.312***
	4	Category (C)	.003	-.002	.175***	.185***	.123***	.001	.001	.001
	5	Distribution Skew (s_k)	.000	.000	.000	-.000	.003	-.293***	-.292***	-.293***
	6	Difficulty (d_f)	-.774***	-.778***	-.566***	-.434***	-.554***	-.389***	-.389***	-.389***
	Effects on Chance Agrt Obsv & Ests	7	Right: Obsd / Estd. Chance Agreement as Dependent Variables Down: Independent Variables	o_{ac}	$ao_{ac}=0^\dagger$	S_{ac}	I_{rac}	AC_{ac}	π_{ac}	κ_{ac}
8		Observed Chance Agreement (o_{ac})	1.00***	---	.021**	.021**	.075***	-.151***	-.152***	-.151***
9		Category (C)	-.019**	---	-.863***	-.863***	-.661***	-.013*	-.014*	-.013*
10		Distribution Skew (s_k)	-.001	---	.000	.000	-.039***	.437***	.434***	.437***
11		Difficulty (d_f)	.585***	---	.000	.000	.009	-.123***	-.125***	-.123***
N	12	N_c (number of rating sessions)	384	384	384	384	384	384	384	384
	13	N_d (number items within each session)	100	100	100	100	100	100	100	100

Main cell entries are directional r squared (dr^2), which are r squared with the directional sign of r , $dr^2=r \cdot |r|$.

*: $p < .05$; **: $p < .01$; ***: $p < .001$. † As ao_{ac} , the chance estimate of a_o , is a constant, its correlations (dr^2) with other variables cannot be calculated.

Interrater Reliability Estimators Tested

Table 4 Mean of Errors (m_e) / Distance Between Index Estimations and Targets of Estimation

		A.	B.	C.	D.	E.	F.	G.	
	1	Author or Source	%- agreement	<i>Bennett et al.</i>	<i>Perreault & Leigh</i>	<i>Gwet</i>	<i>Scott</i>	<i>Cohen</i>	<i>Krippen- dorff</i>
Interrater Reliability	2	Interrater Reliability Estimator	a_o	S	I_r	AC_I	π	κ	α
	3	$m_e(r_i) = \text{mean}(r_i - o_{ri})$ ($0 \leq m_e \leq 1$)	.130***	.096***	.180***	.093***	.327***	.324***	.323***
	4	Standard Deviation of $m_e(r_i)$.145	.099	.148	.104	.221	.220	.220
	5	95% confidence interval of $m_e(r_i)$.115~.144	.086~.106	.164~.194	.082~.103	.304~.349	.302~.346	.301~.345
Chance Agreement	6	Chance Agreement Estimator	a_{oac}	S_{ac}	I_{rac}	AC_{ac}	π_{ac}	κ_{ac}	α_{ac}
	7	$m_e(a_c) = \text{mean}(a_c - o_{ac})$ ($0 \leq m_e \leq 1$)	.130***	.182***	.182***	.130***	.450***	.448***	.448***
	8	Standard Deviation of $m_e(a_c)$.145	.141	.141	.127	.201	.201	.202
	9	95% confidence interval of $m_e(a_c)$.115~.144	.168~.196	.168~.196	.117~.143	.429~.470	.428~.469	.427~.468
N	10	N_c (number of rating sessions)	384	384	384	384	384	384	384
	11	N_d (number items within each session)	100	100	100	100	100	100	100
*: $p < .05$, **: $p < .01$, ***: $p < .001$									

Interrater Reliability Estimators Tested

Table 5 Means and Error of Means (e_m): Index Estimations Against Observations

			A.	B.	C.	D.	E.	F.	G.	H.
	1	Right: Author or Source	<i>Observed Agreement</i>	<i>%-agreement</i>	<i>Bennett et al.</i>	<i>Perreault & Leigh</i>	<i>Gwet</i>	<i>Scott</i>	<i>Cohen</i>	<i>Krippendorff</i>
Interrater Reliability	2	Observed or Estimated Reliability (denotation)	o_{ri}	a_o	S	I_r	AC_i	π	κ	α
	3	Observed / Estimated Interrater Reliability	.555	.685	.556	.726	.600	.237	.240	.241
	4	Standard Deviation	.248	.122	.203	.173	.192	.249	.247	.248
	5	Range (minimum~maximum)	-.20~.90	.42~.92	-.10~.856	.0~.925	-.045~.912	-.177~.778	-.173~.778	-.17~.779
	6	$e_m(r_i)=\text{mean}(r_i)-\text{mean}(o_{ri})$ ($-1 \leq e_m \leq 1$)	.000	.130***	.001	.171***	.044***	-.318***	-.315***	-.314***
	7	95% confidence interval	.00~.00	.115~.144	-.013~.015	.155~.186	.031~.058	-.341~-.295	-.338~-.292	-.338~-.291
	Chance Agreement	8	Chance Agreement (denotation)	o_{ac}	aO_{ac}	S_{ac}	Ir_{ac}	AC_{ac}	π_{ac}	κ_{ac}
9		Observed or Estimated Chance Agreement	.130	.000	.260	.260	.173	.575	.573	.572
10		Standard Deviation	.145	.000	.146	.146	.148	.109	.109	.110
11		Range (minimum~maximum)	.0~.72	.0~.0	.125~.50	.125~.50	.022~.50	.448~.905	.447~.905	.445~.905
12		$e_m(a_c)=\text{mean}(a_c)-\text{mean}(o_{ac})$ ($-1 \leq e_m \leq 1$)	.000	-.130***	.131***	.131***	.044***	.445***	.443***	.443***
13		95% confidence interval	.00~.00	-.144~-.115	.111~.15	.111~.15	.026~.061	.423~.466	.422~.465	.421~.464
N	14	N_c (number of rating sessions)	38	384	384	384	384	384	384	384
	15	N_d (number items within each session)	100	100	100	100	100	100	100	100

*: $p < .05$, **: $p < .01$, ***: $p < .001$

Interrater Reliability Estimators Tested

Table 6 Effects of Category, Skew, and Difficulty on Observed Chance Agreement, Reliability, and Index Estimations (Average Scores)

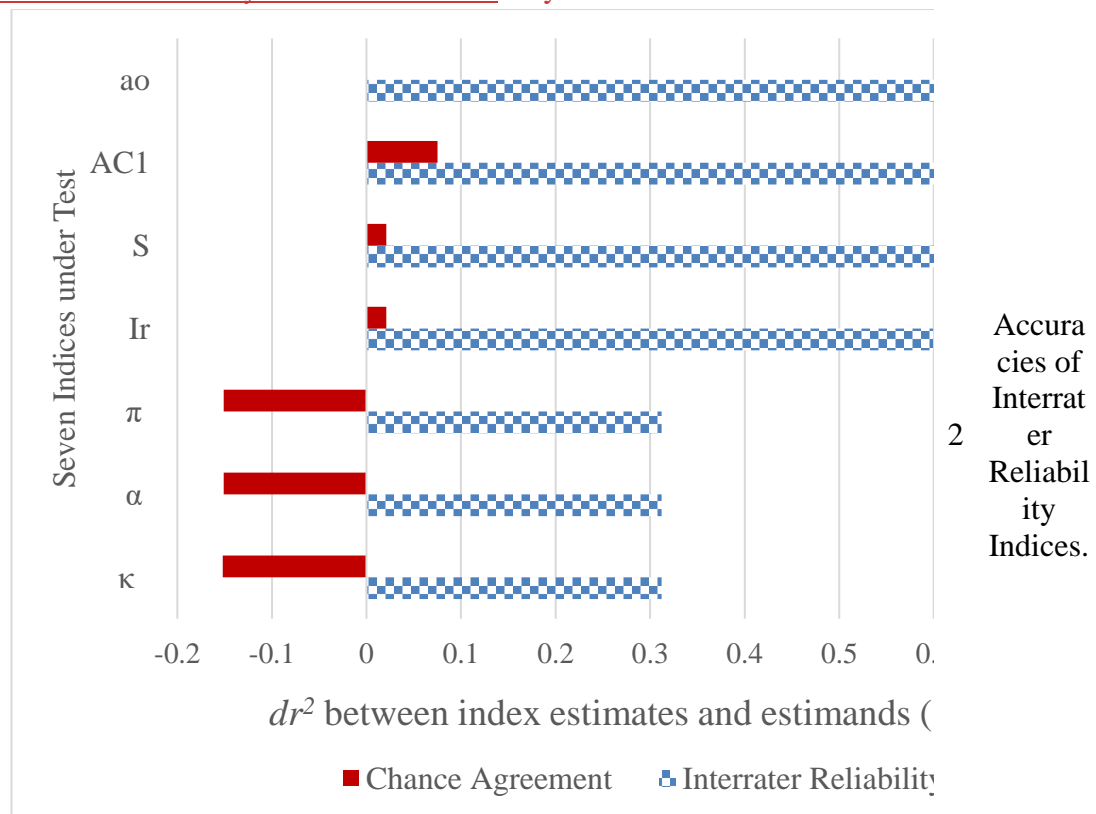
		A.	B.	C.	D.	E.	F.	G.	H.	I.	J.	K.	L.	M.	N.	O.	P.	Q.	
		Reliability Observation or Estimation								Chance Agreement Observation or Estimation									
1	Author/ Source	Observed	%- Agreement	<i>Bennett et al.</i>	<i>Perreault & Leigh</i>	<i>Gwet</i>	<i>Scott</i>	<i>Cohen</i>	<i>Krippen- dorff</i>	Observed	%- Agreement	<i>Bennett et al.</i>	<i>Perreault & Leigh</i>	<i>Gwet</i>	<i>Scott</i>	<i>Cohen</i>	<i>Krippen- dorff</i>		
2	Estimator:	o_{ri}	a_o	S	I_r	AC_1	π	κ	α	o_{ac}	aO_{ac}	S_{ac}	I_{rac}	AC_{ac}	π_{ac}	κ_{ac}	α_{ac}	N_c	
3	Ground 0	.555	.685	.370	.608	.371	.369	.370	.373	.130	0	.500	.500	.499	.501	.500	.498	32	
4	Category (C)	2	.537	.701	.402	.584	.470	.230	.232	.234	.164	0	.500	.500	.401	.598	.597	.596	96
5		4	.550	.678	.571	.747	.621	.226	.230	.230	.128	0	.250	.250	.142	.573	.571	.571	96
6		6	.557	.676	.612	.777	.644	.239	.241	.242	.119	0	.167	.167	.087	.562	.561	.561	96
7		8	.578	.686	.641	.796	.664	.254	.257	.257	.108	0	.125	.125	.062	.564	.563	.562	96
8	Skew (sk)	.50	.550	.688	.560	.732	.592	.370	.372	.374	.138	0	.260	.260	.203	.501	.500	.498	128
9		.75	.556	.678	.547	.722	.588	.302	.304	.305	.122	0	.260	.260	.186	.545	.543	.543	128
10		.99	.560	.690	.561	.723	.619	.040	.044	.045	.130	0	.260	.260	.132	.678	.676	.676	128
11	Difficulty (d_f)	.000	.824	.844	.782	.884	.810	.482	.484	.485	.020	0	.260	.260	.152	.630	.629	.628	48
12		.143	.783	.805	.728	.852	.761	.404	.406	.407	.021	0	.260	.260	.158	.616	.615	.615	48
13		.286	.721	.757	.659	.808	.697	.341	.343	.344	.036	0	.260	.260	.164	.599	.598	.600	48
14		.429	.659	.721	.600	.765	.643	.273	.275	.277	.062	0	.260	.260	.169	.591	.589	.588	48
15		.571	.543	.659	.518	.706	.563	.196	.199	.200	.116	0	.260	.260	.180	.565	.563	.563	48
16		.714	.439	.606	.444	.647	.495	.117	.121	.121	.168	0	.260	.260	.182	.548	.546	.546	48
17		.857	.331	.567	.387	.591	.440	.068	.071	.072	.236	0	.260	.260	.189	.534	.533	.532	48
18		1.00	.142	.523	.332	.552	.389	.018	.022	.022	.380	0	.260	.260	.194	.514	.512	.511	48
19	Mean	.555	.685	.556	.726	.600	.237	.240	.241	.130	0	.260	.260	.173	.575	.573	.572	384	
20	N_d	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	100	

Interrater Reliability Estimators Tested Why the Indices Fail



Figure 1 A sample screen seen by some raters (for category = 6, difficulty = 1).

Interrater Reliability Estimators Tested Why the Indices Fail



Figure

Notes to Figure 2:

1. Solid red bars are dr^2 between estimated chance agreement & observed chance agreement.
2. Dotted blue bars are dr^2 between estimated interrater reliability & observed interrater reliability.
3. Primary benchmark: $dr^2 > 0.8$.
4. Data source: Lines 3 & 8, Table 3.