

1 Agreement Between an Isolated Rater and a Group of
2 Raters

3 S. Vanbelle and A. Albert

4 Medical Informatics and Biostatistics,

5 University of Liège, CHU Sart Tilman (B23),4000 Liège, Belgium

6 Correspondence should be sent to

7 E-Mail: sophie.vanbelle@ulg.ac.be

8 Phone: +3243662590

9 Fax: +3243662596

10

11

Abstract

12

13

14

15

16

17

18

19

20

The agreement between two raters judging items on a categorical scale is traditionally measured by Cohen's kappa coefficient. We introduce a new coefficient for quantifying the degree of agreement between an isolated rater and a group of raters on a nominal or ordinal scale. The coefficient, which is defined on a population-based model, requires a specific definition of the concept of perfect agreement but possesses the same properties as Cohen's kappa coefficient. Further, it reduces to the classical kappa when there is only one rater in the group. An intraclass and a weighted versions of the coefficient are also introduced. The new approach overcomes the problem of

21 consensus and generalizes Schouten's index. The sampling variability of the
22 agreement coefficient is derived by the Jackknife technique. The method is
23 illustrated on published syphilis data and on data collected from a study
24 assessing the ability of medical students in diagnostic reasoning.

25 Keywords: kappa coefficient; nominal scale; ordinal scale.

1 INTRODUCTION

26

27 Cohen (1960) introduced the kappa coefficient $\kappa = (p_o - p_e)/(1 - p_e)$ to quantify
28 the agreement between two raters classifying N items on a binary or nominal
29 scale. He corrected the proportion of items with concordant classification (p_o)
30 for the proportion of concordant pairs expected by chance (p_e) and standardized
31 the quantity to obtain a value 1 when the agreement between the two raters is
32 perfect and 0 when the observed agreement is equal to the agreement expected
33 by chance. There are situations where the agreement between an isolated rater
34 and a group of raters is needed. For example, each of a series of individuals may
35 be assessed against a group of experts and a ranking of the individuals may be
36 required. Conversely, agreement may be searched between a group of users and a
37 gold standard. Usually in such instances, a consensus is determined in the group of
38 raters and the problem is reduced to the case of measuring agreement between the
39 isolated rater and the consensus in the group (Landis and Koch 1977, Soeken and
40 Prescott 1986, Salerno *et al.* 2003). The consensus may be defined as the category
41 chosen by a given proportion of the raters in the group (for example, Ruperto *et*
42 *al.* (2006) defined the consensus as the category chosen by at least 80% of the
43 raters in the group) or the category the most frequently chosen by the raters in
44 the group (Kalant *et al.* (2000), Smith *et al.* (2003)). In both cases, the problem of
45 how to handle items without consensus arises. Ruperto *et al.* (2006) discarded all
46 patients without consensus from the analysis, while Kalant *et al.* (2000) and Smith

47 *et al.* (2003) did not encounter the problem. The method consisting in reducing
48 the judgements made by a group of raters into a consensus decision was criticized
49 by Eckstein *et al.* (1998), Salerno *et al.* (2003) and Miller *et al.* (2004). Eckstein
50 *et al.* (1998) investigated the bias that may result from removing items without
51 consensus, while Salerno *et al.* (2003) argued that the dispersion likely to occur
52 in the classifications made by the raters in the group may not be reflected in the
53 consensus. Finally, Miller *et al.* (2004) examined the possibility to obtain different
54 conclusions by using different rules of consensus. Light (1971) developed a statistic
55 for comparing the joint agreement of several raters with a gold standard. This
56 statistic is a mixture of the proportions of concordant pairs obtained between each
57 of the rater in the group and the gold standard (the isolated rater). His method
58 leads to tedious calculations, does not quantify the agreement between the gold
59 standard and the group of raters and the calculations have not been extended to
60 the case of a group including more than 3 raters. Williams (1976) developed a
61 measure for comparing the joint agreement of several raters with another rater
62 without determining a consensus in the group of raters. Indeed, he compared the
63 mean proportion of concordant items between the isolated rater and each rater
64 in the group to the mean proportion of concordant items between all possible
65 pairs of raters among the group of raters. The ratio derived (Williams' index) is
66 compared to the value of 1. Unfortunately, the coefficient proposed by Williams
67 (1976) does not correct for agreements due to chance and does not quantify the

68 agreement between the isolated rater and the group of raters. Finally, Schouten
69 (1982) developed a method of hierarchical clustering based on pairwise weighted
70 agreement measures to select one or more homogeneous subgroups of raters when
71 several raters classify items on a nominal or an ordinal scale. Hereafter, we propose
72 a coefficient for quantifying the agreement between an isolated rater and a group
73 of raters, which overcomes the problem of consensus, generalizes the approach of
74 Schouten (1982) and possesses the same properties as Cohen's kappa coefficient.

75 **2 DEFINITION OF THE AGREEMENT INDEX**

76 **2.1 Binary scale (K=2)**

77 Consider a population \mathcal{I} of items and a population \mathcal{R} of raters. Suppose that the
78 items have to be classified on a binary scale by the population of raters and by an
79 independent isolated rater. Let $X_{i,r}$ be the random variable such that $X_{i,r} = 1$ if a
80 randomly selected rater r of the population \mathcal{R} classifies a randomly selected item
81 i of population \mathcal{I} in category 1 and $X_{i,r} = 0$ otherwise. Let $E(X_{i,r}) = P(X_{i,r} =$
82 $1) = p_i$ over the population of raters. Then, over the population of items, let
83 $E(p_i) = \pi$ and $\sigma^2 = var(p_i)$. In the same way, let Y_i denote the random variable
84 equal to 1 if the isolated rater classifies item i in category 1 and $Y_i = 0$ otherwise.
85 Over the population of items, $E(Y_i) = \pi^*$ and $var(Y_i) = \sigma^{*2} = \pi^*(1 - \pi^*)$. Finally,
86 let ICC denote the intraclass correlation coefficient in the population of raters

Table 1: Theoretical model for the classification of a randomly selected item i on a binary scale by the population of raters \mathcal{R} and the isolated rater

		Isolated rater		
\mathcal{R}		0	1	
0	$E[(1 - p_i)(1 - Y_i)]$	$E[(1 - p_i)Y_i]$	$1 - \pi$	
	$(1 - \pi)(1 - \pi^*) + \rho\sigma\sigma^*$	$(1 - \pi)\pi^* - \rho\sigma\sigma^*$		
1	$E[p_i(1 - Y_i)]$	$E[p_iY_i]$	π	
	$\pi(1 - \pi^*) - \rho\sigma\sigma^*$	$\pi\pi^* + \rho\sigma\sigma^*$		
		$1 - \pi^*$	π^*	1

87 (Fleiss 1981)

$$ICC = \frac{\sigma^2}{\pi(1 - \pi)} \quad (1)$$

88 and ρ the correlation between p_i and Y_i over \mathcal{I}

$$\rho = \frac{E(p_iY_i) - \pi\pi^*}{\sigma\sigma^*}. \quad (2)$$

89 Using these definitions, a 2×2 table can be constructed cross-classifying the popu-
 90 lation of raters \mathcal{R} and the isolated rater with respect to the binary scale (Table 1).

91

92 The probability that the population of raters and the isolated rater agree on
 93 item i is defined by

$$\Pi_i = p_iY_i + (1 - p_i)(1 - Y_i) \quad (3)$$

94 so that, over the population of items \mathcal{I} , the mean probability of agreement is given
 95 by the expression

$$\Pi_T = E(\Pi_i) = \pi\pi^* + (1 - \pi)(1 - \pi^*) + 2\rho\sigma\sigma^* \quad (4)$$

96 By definition, the population of raters and the isolated rater are considered to be
 97 in "perfect agreement" if and only if

$$\pi = \pi^* = \pi^{**} \text{ and } \rho = 1. \quad (5)$$

98 In terms of the random variables p_i and Y_i over \mathcal{I} this is equivalent to writing

$$p_i = \pi^{**}(1 - \sqrt{ICC}) + \sqrt{ICC}Y_i \quad (6)$$

99 It follows from Equation 4 that the maximum attainable probability of perfect
 100 agreement is given by

$$\Pi_M = 1 - 2\pi^{**}(1 - \pi^{**})(1 - \sqrt{ICC}) \quad (7)$$

101 which turns out to be equal to 1 only if $ICC = 1$, i.e. that there is perfect
 102 agreement between all raters in population \mathcal{R} , or trivially if $\pi^{**} = 0$ or 1.

103 Then, the coefficient of agreement between the population of raters and the
 104 isolated rater is defined in a kappa-like way:

$$\kappa = \frac{\Pi_T - \Pi_E}{\Pi_M - \Pi_E} \quad (8)$$

105 where Π_E is the agreement expected by chance, i.e., the probability that the pop-
 106 ulation of raters and the isolated rater agree under the independence assumption

107 $(E(p_i Y_i) = E(p_i)E(Y_i))$, defined by

$$\Pi_E = \pi\pi^* + (1 - \pi)(1 - \pi^*) \quad (9)$$

108 Note that $\Pi_T = \Pi_E$ when there is no correlation between the ratings of the pop-
 109 ulation of raters and the isolated rater ($\rho = 0$) or when there is no variability in
 110 the classification made by the populations of raters ($\sigma^2 = 0$) or by the isolated
 111 rater ($\sigma^{*2} = 0$).

112

113 An intraclass version of the agreement index κ_I may be derived by assuming
 114 that $\pi = \pi^* = \pi^{**}$. It leads to

$$\kappa_I = \rho = \frac{E(p_i Y_i) - \pi^{**}}{\sigma\sqrt{\pi^{**}(1 - \pi^{**})}} \quad (10)$$

115 **2.2 Multinomial scale (K>2)**

116 When $K > 2$, the coefficient of agreement between the population of raters and
 117 the isolated rater is defined by

$$\kappa = \frac{\sum_{j=1}^K (\Pi_{[j]T} - \Pi_{[j]E})}{\sum_{j=1}^K (\Pi_{[j]M} - \Pi_{[j]E})} = \frac{\Pi_T - \Pi_E}{\Pi_M - \Pi_E} \quad (11)$$

118 where $\Pi_{[j]T}$, $\Pi_{[j]E}$ and $\Pi_{[j]M}$ correspond to the quantities described in the 2×2
 119 case when the nominal scale is dichotomized by grouping all categories other than
 120 category j together and Π_T , Π_E and Π_M are defined respectively by

$$\Pi_T = \sum_{j=1}^K E(p_{ij} Y_{ij}); \quad \Pi_E = \sum_{j=1}^K \pi_j \pi_j^*;$$

$$\Pi_M = \sum_{j=1}^K E((\pi_j^{**} + (1 - \pi_j^{**})\sqrt{ICC_j})Y_{ij}) \quad (12)$$

121 where p_{ij} denotes the probability for a randomly selected item i to be classified
 122 in category j ($j = 1, \dots, K$) by the population of raters with $E(p_{ij}) = \pi_j$ and
 123 Y_{ij} denotes the random variable equal to 1 if the isolated rater classifies item i in
 124 category j ($Y_{ij} = 0$ otherwise). Finally, ICC_j denotes the intraclass correlation
 125 coefficient relative to category j ($j = 1, \dots, K$) in the population of raters.

126

127 The coefficient κ possesses the same properties as Cohen's kappa coefficient,
 128 $\kappa = 1$ when agreement is perfect ($\Pi_T = \Pi_M$), $\kappa = 0$ if observed agreement is equal
 129 to agreement expected by chance ($\Pi_T = \Pi_E$) and $\kappa < 0$ if observed agreement is
 130 lower than expected by chance ($\Pi_T < \Pi_E$).

131 2.3 Ordinal scale ($K > 2$)

132 A weighted version of the agreement index can be defined in a way similar to the
 133 weighted kappa coefficient (Cohen 1968),

$$\kappa_W = \frac{\Pi_{T,W} - \Pi_{E,W}}{\Pi_{M,W} - \Pi_{E,W}} \quad (13)$$

134 with

$$\Pi_{T,W} = \sum_{j=1}^K \sum_{k=1}^K w_{jk} E(p_{ij} Y_{ik}); \quad (14)$$

135

$$\Pi_{E,W} = \sum_{j=1}^K \sum_{k=1}^K w_{jk} \pi_j \pi_k^*; \quad (15)$$

136

$$\Pi_{M,W} = \sum_{j=1}^K \sum_{k=1}^K w_{jk} E(\pi^{**}(1 - \sqrt{ICC_j})Y_{ik}) + \sum_{k=1}^K \sqrt{ICC_k}. \quad (16)$$

137

138

139

In general, $0 \leq w_{jk} \leq 1$ and $w_{kk} = 1$, ($j, k = 1, \dots, K$). Cicchetti and Allison (1971) have defined absolute weights $w_{jk} = 1 - \frac{|j-k|}{K-1}$ whereas Fleiss and Cohen (1973) suggested quadratic weights $w_{jk} = 1 - \left(\frac{j-k}{K-1}\right)^2$.

140

3 ESTIMATION OF THE PARAMETERS

141

142

143

Suppose that a random sample of N items drawn from population \mathcal{I} is classified on a K -categorical scale by a random sample (group) of size R from the population of raters \mathcal{R} and by an independent isolated rater.

144

3.1 Binary scale ($K = 2$)

145

146

147

148

149

150

151

152

Let $x_{i,r}$ denotes the observed value of the random variable $X_{i,r}$ denoting the rating of rater r of the population \mathcal{R} ($i = 1, \dots, N; r = 1, \dots, R$). Let y_i denotes the observed value of the random variable Y_i representing the rating of the isolated rater. Then, let $n_i = \sum_{r=1}^R x_{i,r}$ denotes the number of times the item i is classified in category 1 by the group of raters and let $\hat{p}_i = n_i/R$ denote the corresponding proportions ($i = 1, \dots, N$).

The intraclass correlation coefficient in the group of raters is estimated by

153 (Fleiss 1981)

$$\widehat{ICC} = 1 - \frac{\sum_{i=1}^N n_i(R - n_i)}{RN(N - 1)p(1 - p)} \quad (17)$$

where p is the proportion of items classified in category 1 by the group of raters,

$$p = \frac{1}{N} \sum_{i=1}^N \hat{p}_i.$$

154 The probability that the population of raters and the isolated rater agree is
155 estimated by the *observed proportion of agreement*,

$$\hat{\Pi}_T = p_o = \frac{1}{N} \sum_{i=1}^N (\hat{p}_i y_i + (1 - \hat{p}_i)(1 - y_i)). \quad (18)$$

156 Clearly, $p_o = 1$ if the raters of the group and the isolated rater classify each item
157 in the same category and $p_o = 0$ if the isolated rater systematically classifies items
158 in a category never chosen by the group of raters.

159

160 The probability of agreement expected by chance is estimated by the *propor-*
161 *tion of agreement expected by chance*,

$$p_e = py + (1 - p)(1 - y) \quad (19)$$

where y is the proportion of items classified in category 1 by the isolated rater,

$$y = \frac{1}{N} \sum_{i=1}^N y_i.$$

162 The degree of agreement κ between the group of raters and the isolated rater
163 is then estimated by

$$\hat{\kappa} = \frac{p_o - p_e}{p_m - p_e} \quad (20)$$

164 where p_m corresponds to the maximum possible proportion of agreement derived
 165 by the data. We have

$$p_m = \frac{1}{N} \sum_{i=1}^N \max(\hat{p}_i, 1 - \hat{p}_i). \quad (21)$$

166 3.2 Multinomial scale ($K > 2$)

167 The estimation of the parameters easily extends to the case $K > 2$. Let $x_{ij,r}$ denote
 168 the observed value of the random variable $X_{ij,r}$ equal to 1 if rater r ($r = 1, \dots, R$)
 169 of the group classified item i ($i = 1, \dots, N$) in category j ($j = 1, \dots, K$) and equal
 170 to 0 otherwise. In the same way, let y_{ij} denote the observed value of the random
 171 variable Y_{ij} corresponding to the rating of the isolated rater. Let $n_{ij} = \sum_{r=1}^R x_{ij,r}$
 172 denotes the number of times the item i is classified in category j by the raters of
 173 the group and let \hat{p}_{ij} denote the corresponding proportions. We have $\sum_{j=1}^K \hat{p}_{ij} = 1$,
 174 ($i = 1, \dots, N$). The data can be conveniently summarized in a 2-way classification
 175 table (see Table 2) by defining the quantities

$$c_{jk} = \frac{1}{N} \sum_{i=1}^N \hat{p}_{ij} y_{ik}, \quad j, k = 1, \dots, K \quad (22)$$

176 The *observed proportion of agreement* between the group of raters and the
 177 isolated rater is defined by

$$p_o = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K \hat{p}_{ij} y_{ij} = \sum_{j=1}^K c_{jj} \quad (23)$$

Table 2: Two-way classification table of the N items by the group of raters and the isolated rater

Group of raters	Isolated rater					Total
	1	...	j	...	K	
1	c_{11}	...	c_{1j}	...	c_{1K}	$c_{1.}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
j	c_{j1}	...	c_{jj}	...	c_{jK}	$c_{j.}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
K	c_{K1}	...	c_{Kj}	...	c_{KK}	$c_{K.}$
Total	$c_{.1}$...	$c_{.j}$...	$c_{.K}$	1

178 The marginal classification distribution of the isolated rater, namely,

$$y_j = \frac{1}{N} \sum_{i=1}^N y_{ij}, \quad j = 1, \dots, K \quad (24)$$

179 with $\sum_{j=1}^K y_j = 1$ and the marginal classification distribution of the group of raters,

$$p_j = \frac{1}{N} \sum_{i=1}^N \hat{p}_{ij}, \quad j = 1, \dots, K \quad (25)$$

180 with $\sum_{j=1}^K p_j = 1$ are needed to estimate the agreement expected by chance. The

181 *proportion of agreement expected by chance* is given by

$$p_e = \sum_{j=1}^K p_j y_j = \sum_{j=1}^K c_{j.} c_{.j} \quad (26)$$

182 The degree of agreement κ between the population of raters and the isolated

183 rater is then estimated by

$$\hat{\kappa} = \frac{p_o - p_e}{p_m - p_e} \quad (27)$$

184 where p_m corresponds to the maximum possible proportion of agreement derived
 185 from the data,

$$p_m = \frac{1}{N} \sum_{i=1}^N \max_j p_{ij}. \quad (28)$$

186 Note that when $R = 1$, $p_m = 1$ and the agreement coefficient $\hat{\kappa}$ reduces to the
 187 classical Cohen's kappa coefficient defined in the case of two isolated raters.

188

189 The intraclass correlation coefficient in the group of raters is estimated by
 190 (Fleiss 1981)

$$\widehat{ICC} = 1 - \frac{NR^2 - \sum_{i=1}^N \sum_{j=1}^K n_{ij}^2}{NR(R-1) \sum_{j=1}^K p_j(1-p_j)} \quad (29)$$

191 3.3 Ordinal scale ($K > 2$)

192 The estimation of the weighted agreement index is done by merely introducing
 193 weights in the estimations previously defined. Hence,

$$\hat{\kappa}_W = \frac{p_{o,w} - p_{e,w}}{p_{m,w} - p_{e,w}} \quad (30)$$

194 with

$$\begin{aligned} p_{o,w} &= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K \sum_{k=1}^K w_{jk} p_{ij} y_{ik} \\ p_{e,w} &= \sum_{j=1}^K \sum_{k=1}^K w_{jk} p_j y_k \\ p_{m,w} &= \frac{1}{N} \sum_{i=1}^N \max_j \left(\sum_{k=1}^K w_{jk} p_{ik} \right). \end{aligned} \quad (31)$$

195 The unweighted agreement index $\hat{\kappa}$ can be obtained using the weights $w_{jj} = 1$
 196 and $w_{jk} = 0, j \neq k$.

197 4 ASYMPTOTIC SAMPLING VARIANCE

198 The Jackknife method (Efron and Tibshirani, 1993) was used to determine the
 199 sampling variance of the agreement index. Suppose that the agreement between
 200 the isolated rater and the population of raters was estimated on a random sample
 201 of N items. Let $\hat{\kappa}_N$ denote that agreement index and $\hat{\kappa}_{N-1}^{(i)}$ denote the estimated
 202 agreement index when observation i is deleted. These quantities are used to de-
 203 termine the pseudo-values

$$\hat{\kappa}_{N,i} = N\hat{\kappa}_N - (N-1)\hat{\kappa}_{N-1}^{(i)} \quad (32)$$

204 The Jackknife estimator of the agreement index is then defined by

$$\tilde{\kappa}_N = \frac{1}{N} \sum_{i=1}^N \hat{\kappa}_{N,i} \quad (33)$$

205 with variance

$$var(\tilde{\kappa}_N) = \frac{1}{N} \left\{ \frac{1}{N-1} \sum_{i=1}^N (\hat{\kappa}_{N,i} - \hat{\kappa}_N)^2 \right\} \quad (34)$$

206 The bias of the Jackknife estimator is estimated by

$$Bias(\tilde{\kappa}_N) = (N-1) \{ \tilde{\kappa}_N - \hat{\kappa}_N \} \quad (35)$$

207 5 CONSENSUS APPROACH

208 The consensus approach consists in summarizing the responses given by the raters
209 of the group in a unique quantity. Most approaches define the modal category (ma-
210 jority rule) or the category chosen by a prespecified proportion of raters ($\geq 50\%$)
211 as the consensus category. A random variable Z_{ij} is then defined to be equal to
212 1 if category j corresponds to the consensus category given by the population \mathcal{R}
213 of raters for item i and equal to 0 otherwise. It is obvious that a consensus may
214 not always be defined. For example, on a multinomial scale, we could have two
215 modal categories or no category chosen by the prespecified proportion of raters.
216 Therefore, suppose that on the N items drawn from population \mathcal{I} , a consensus can
217 only be defined on $N_C \leq N$ items. Let \mathcal{I}_C denote the sub-population of items on
218 which a consensus is always possible. If z_{ij} denotes the observed value of the ran-
219 dom variable Z_{ij} , we have $\sum_{j=1}^K z_{ij} = 1$ and the agreement between the population
220 of raters and the isolated raters is reduced to the case of 2 isolated raters. The
221 Cohen intraclass or weighted kappa coefficient can then be estimated. Note that
222 the strenght of the consensus is not taken into account by the random variable
223 Z_{ij} . For example on a binary scale, using the majority rule, we will have $Z_{ij} = 1$
224 if $p_{ij} = 0.6$ but also if $p_{ij} = 0.9$. It can easily be shown that the new method-
225 ology defined and the consensus approach are equivalent only in two particular
226 cases, firstly when there is only one rater in the group of raters ($R = 1$) and
227 secondly when $\mathcal{I}_C = \mathcal{I}$ and there is perfect agreement in the population of raters

228 ($ICC = 1$).

229

6 EXAMPLES

230 6.1 Syphilis serology

231 A proficiency testing program for syphilis serology was conducted by the College
232 of American Pathologists (CAP). For the fluorescent treponemal antibody absorp-
233 tion test (FTA-ABS), 3 reference laboratories were identified and considered as
234 experts in the use of that test. During 1974, 40 syphilis serology specimens were
235 tested independently by the 3 reference laboratories. Williams (1976) presented
236 data for 28 specimens. To evaluate the performance of a participant, the agree-
237 ment between the participant and the 3 reference laboratories had to be evaluated.
238 The data are summarized in a two-way classification table (Table 3) as explained
239 in section 2.3.

240 Using the quadratic weighting scheme, the weighted coefficient of agreement
241 $\hat{\kappa}_W$ amounted 0.79 ± 0.06 . When applying the consensus approach based on the
242 majority rule, we found a weighted kappa coefficient of 0.76 ± 0.06 . Remark that
243 2 specimens were eliminated because no consensus was found in the group of the
244 3 reference laboratories. Finally, the weighted agreement measure developed by
245 Schouten (1982) was 0.73 ± 0.07 . Note that the intraclass correlation coefficient
246 was 0.68 ± 0.06 in the group of raters.

Table 3: Two-way classification table of the 28 syphilis serology specimens as NR (non-reactive), B (borderline) and R (reactive) by 3 reference laboratories and a participant

	Participant			
Reference laboratories	NR	B	R	Total
NR	0.143	0.250	0.024	0.417
B	0	0.036	0.071	0.107
R	0	0	0.476	0.476
Total	0.143	0.286	0.571	1

247 6.2 Script Test of Concordance

248 The Script Test of Concordance (SCT) is used in medicine to evaluate the ability
249 of physicians or medical students (isolated raters) to solve clinical situations not
250 clearly defined (Charlin *et al.* 2002). The complete test consists of a number of
251 items $(1, \dots, N)$ to be evaluated on a 5-point Likert scale ($K = 5$). Each item
252 represents a clinical situation likely to be seen in real life practice and a poten-
253 tial assumption is proposed with it. The situation has to be unclear, even for an
254 expert. The task of the student or the physician being evaluated is to consider
255 the effect of additional evidence on the suggested assumption. In this respect, the
256 candidate has to choose between the following proposals: (-2) The assumption is
257 practically eliminated; (-1) The assumption becomes less likely; (0) The informa-

258 tion has no effect on the assumption; (+1) The assumption becomes more likely
259 and (+2) The assumption is practically the only possibility. The questionnaire is
260 also given to a panel of experts (raters $1, \dots, R$). The problem is to evaluate the
261 agreement between each individual medical student and the panel of experts.

262

263 Between 2003 and 2005, the SCT was proposed to students specializing in gen-
264 eral practice at the University of Liège, Belgium (Vanbelle *et al.* 2007). The SCT
265 consisted of 34 items relating possible situations encountered in general practice.
266 There were 39 students passing the test and completing the entire questionnaire.
267 Their responses were confronted to the responses of a panel of 11 experts. The
268 intraclass correlation coefficient was 0.22 ± 0.04 in the group of experts. Using
269 the quadratic weighting scheme, the individual $\hat{\kappa}_W$ coefficients for the 39 students
270 ranged between 0.37 and 0.84. The mean value ($\pm SD$) was 0.61 ± 0.12 .

271 Using the consensus method, where consensus was defined as either the major-
272 ity of the raters or a proportion of at least 50% of the raters, respectively 2 (6%)
273 and 12 (35%) items had to be omitted from the analysis because no consensus was
274 reached among the raters. The mean weighted kappa values for the 39 students
275 were equal to 0.49 ± 0.13 (range: 0.19-0.72) and 0.66 ± 0.14 (range: 0.23-0.82)
276 with the majority and the 50% rules, respectively. Figure 1 displays the individ-
277 ual agreement coefficients relative to each student for the different methods. A
278 ranking of the student was needed in order to select only the best students. The

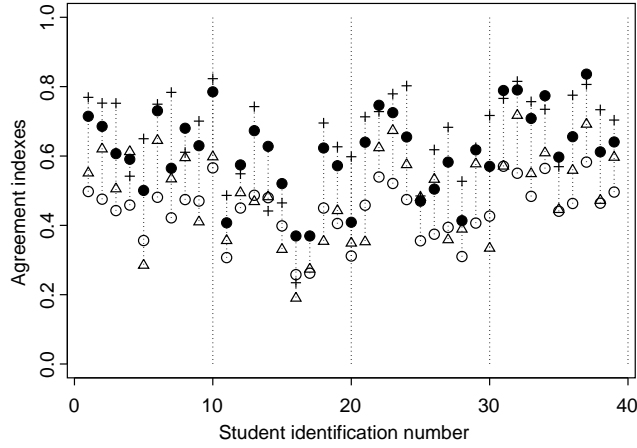


Figure 1: Values of κ_W (●), weighted κ coefficients using the majority (Δ) and the 50% (+) rules and weighed agreement index of Schouten (○) for the 39 students passing the SCT

279 ranking changed markedly for some students according to the method used. For
 280 example, student No. 39 ranked at the 16th place with the new approach, the 9th
 281 place with Schouten index, the 10th place using the majority rule and at 20th
 282 place using the 50% rule.

283 7 DISCUSSION

284 The method described in this paper was developed to quantify the agreement be-
 285 tween an isolated rater and a group of raters judging items on a categorical scale.
 286 A population-based approach was used but in case of a fixed group of raters, es-
 287 timates are replaced by actual values. The derived agreement index κ possesses

288 the same properties as Cohen's kappa coefficient (Cohen 1960) and reduces to it
289 if there is only one rater in the group. The isolated rater and the group of raters
290 are defined to be in perfect agreement when they have the same probability, for
291 each item, to classify this item in a given category and the correlation coefficient
292 between the isolated rater and the population of raters is equal to 1. It can be
293 shown that with the additional assumption of perfect agreement in the population
294 of raters ($ICC = 1$), the proposed agreement index κ is algebraically equivalent
295 to the agreement coefficient derived by Schouten (1982). In other terms, perfect
296 agreement can be reached between the isolated rater and the population of raters
297 even if no perfect agreement occurs in the population of raters unlike the agree-
298 ment index of Schouten (1982). The new approach is equivalent to the consensus
299 approach when it is possible to determine a consensus for all items of the sample
300 and there is perfect agreement in the group of raters on each item. The proposed
301 method is superior the consensus approach in the sense that no decision has to
302 be made if there is no consensus in the group. Moreover, the new approach takes
303 into account the variability in the group while the strength of consensus is not
304 taken into account with the consensus method. Finally, as illustrated in the SCT
305 example and pointed out by Salerno *et al.* (2003) and Miller *et al.* (2004), the re-
306 sults may vary substantially according to the definition of the consensus used. The
307 proposed kappa coefficient thus provides an alternative to the common approach
308 which consists in summarizing the responses given by the raters in the group into

309 a single response (the consensus) and generalizes the agreement index proposed
310 by Schouten (1982). Further, it has the advantage of using more information than
311 the consensus method (variability in the group of raters), of solving the problem
312 of items without consensus and of being built upon less stringent assumptions.
313 Experts can fix levels to interpret the values taken by the new coefficient and
314 determine a lower bound under which the isolated rater may be rejected as in the
315 SCT selection process or considered as "out of range".

316 References

317 CHARLIN, B., R. GAGNON, L. SIBERT AND C. VAN DER VELUTEN (2002), Le
318 test de concordance de script : un instrument d'évaluation du raisonnement
319 clinique, *Pédagogie Médicale* 3, 135–144.

320 CICHETTI, D.V. AND T. ALLISON (1971), A new procedure for assessing reli-
321 ability of scoring EEG sleep recordings, *American Journal of EEG Technology*
322 11, 101–109.

323 COHEN, J. (1960), A coefficient of agreement for nominal scales, *Educational and*
324 *Psychological Measurement* 20, 37–46.

325 COHEN, J. (1968), Weighted Kappa: nominal scale agreement with provision for
326 scaled disagreement or partial credit, *Psychological Bulletin* 70, 213–220.

327 ECKSTEIN, M.P., T.D. WICKENS, G. AHARONOV, G. RUAN, C.A. MORIOKA
328 AND J.S. WHITING (1998), Quantifying the limitations of the use of consensus
329 expert committees in ROC studies, *Proceedings SPIE: Medical Imaging 1998:*
330 *Image perception* 3340, 128–134.

331 EFRON, B. AND R.J. TIBSHIRANI (1993), *An introduction to the bootstrap*, Chap-
332 man and Hall, New York.

333 FLEISS, J.L. AND J. COHEN (1973), The equivalence of weighted kappa and
334 the intraclass correlation coefficient as measure of reliability, *Educational and*
335 *psychological measurement* 33, 613–619.

336 FLEISS, J.L. (1981), *Statistical methods for rates and proportions* (2nd edition),
337 John Wiley, New York

338 KALANT, N., M. BERLINGUET, J.G. DIODATI, L. DRAGATAKIS AND F. MAR-
339 COTTE (2000), How valid are utilization review tools in assessing appropriate
340 use of acute care beds? *Canadian Medical Association Journal* 162, 1809–1813.

341 LANDIS, J.R. AND G.G. KOCH (1977), An application of hierarchical kappa-type
342 statistics in the assessment of majority agreement among multiple observers,
343 *Biometrics* 33, 363–374.

344 LIGHT, R.J. (1971), Measures of response agreement for qualitative data: some
345 generalizations and alternatives, *Psychological bulletin* 76, 365–377.

346 MILLER, D.P., K.F. O'SHAUGHNESSY, S.A. WOOD AND R.A. CASTELLINO
347 (2004), Gold standards and expert panels: a pulmonary nodule case study with
348 challenges and solutions, *Proceedings SPIE: Medical Imaging 1998: Image per-*
349 *ception, Observer Performance and Technology Assessment* 5372, 173–184.

350 RUPERTO, N., A. RAVELLI, S. OLIVEIRA, M. ALESSIO, D. MIHAYLOVA,
351 S. PASIC, E. CORTIS, M. APAZ, R. BURGOS-VARGAS, F. KANAKOUDI-
352 TSAKALIDOU, X. NORAMBUENA, F. CORONA, V. GERLONI, S. HAGELBERG,
353 A. AGGARWAL, P. DOLEZALOVA, C.M. SAAD, S.C. BAE, R. VESELY, T.
354 AVCIN, H. FOSTER, C. DUARTE, T. HERLIN, G. HORNEFF, L. LEPORE,
355 M. VAN ROSSUM, L. TRAIL, A. PISTORIO, B. ANDERSSON-GARE, E.H. GI-
356 ANNINI AND A. MARTINI (2006), Pediatric Rheumatology International Tri-
357 als Organization. The pediatric Rheumatology International Trials Organiza-
358 tion/American College of Response to Therapy in Juvenile Systemic Lupus Ery-
359 thematosus: Prospective Validation of the Definition of Improvement, *Arthritis*
360 *and Rheumatism (Arthritis Care and Research)* 55, 355–363.

361 SALERNO, S.M., P.C. ALGUIRE AND S.W. WAXMAN (2003), Competency in
362 interpretation of 12-Lead Electrocardiograms: A summary and Appraisal of
363 Published Evidence, *Annals of Internal Medicine* 138, 751–760.

364 SCHOUTEN, H.J.A. (1982), Measuring pairwise interobserver agreement when all
365 subjects are judged by the same observers, *Statistica Neerlandica* 36, 45–61.

- 366 SMITH, R., A.J. COPAS, M. PRINCE, B. GEORGE, A.S. WALKER AND S.T.
367 SADIQ (2003), Poor sensitivity and consistency of microscopy in the diagnosis
368 of low grade non-gonococcal urethritis, *Sexually Transmitted Infections* 79, 487–
369 490.
- 370 SOEKEN, K.L. AND P.A. PRESCOTT (1986), Issues in the use of kappa to estimate
371 reliability, *Medical care* 24, 733–741.
- 372 VANBELLE, S., V. MASSART, D. GIET AND A. ALBERT (2007), Test de concor-
373 dance de script: un nouveau mode d'établissement des scores limitant l'effet du
374 hasard, *Pédagogie Médicale* 8, 71–81.
- 375 WILLIAMS, G.W. (1976), Comparing the Joint Agreement of Several Raters with
376 Another Rater, *Biometrics* 32, 619–627.