

1 **A note on the linearly weighted kappa**
2 **coefficient for ordinal scales**

3 S. Vanbelle and A. Albert

4 *Medical Informatics and Biostatistics, School of Public Health, University of*
5 *Liège, Liège, Belgium*

6 **Abstract**

7 A frequent criticism formulated against the use of weighted kappa coefficients is
8 that the weights are arbitrarily defined. We show that using linear weights for a K-
9 ordinal scale is equivalent to deriving a kappa coefficient from K-1 embedded 2×2
10 tables.

11 *Key words:* absolute weights, interpretation, agreement, disagreement

12 **1 INTRODUCTION**

13 Cohen's kappa coefficient (Cohen, 1960) is widely used to quantify agreement
14 between two raters on a nominal scale (Ludbrook, 2002). It corrects the ob-
15 served percentage of agreements between the raters for the effect of chance.

* Sophie Vanbelle, Medical Informatics and Biostatistics, School of Public Health, University of Liège, CHU Sart Tilman, 4000 Liège, Belgium, sophie.vanbelle@ulg.ac.be

16 A value of 0 implies no agreement beyond chance, whereas a value of 1 cor-
17 responds to a perfect agreement between the two raters. There are situations
18 where disagreements between raters may not all be equally important. For
19 example, on an ordinal scale, a greater "penalty" will be applied if the two
20 categories chosen by the raters are farther apart. To account for these inequal-
21 ities, Cohen (1968) introduced weights in the formulation of the agreement
22 index leading to the weighted kappa coefficient. Although the weights are in
23 general arbitrarily defined, those introduced by Cicchetti and Allison (1971)
24 and by Fleiss and Cohen (1973) are the most commonly used. The former
25 are linear and the latter have a quadratic form. Cohen (1968) showed that,
26 under specific conditions, the weighted kappa coefficient is equivalent to the
27 product-moment correlation coefficient. Moreover, Fleiss and Cohen (1973)
28 and Schuster (2004) showed that the weighted kappa with a quadratic weight-
29 ing scheme is equivalent to the intraclass correlation coefficient. Hereafter,
30 we show that the weighted kappa coefficient defined with linear weights for a
31 K-ordinal scale can be derived from (K-1) embedded 2×2 contingency tables.

32 2 DEFINITION OF THE WEIGHTED KAPPA COEFFICIENT

33 Consider two raters who classify a sample of n subjects (or objects) into K
34 categories of an ordinal scale (see Table 1), where n_{ij} is the number of items
35 classified into category i by rater 1 and category j by rater 2, $n_{i\cdot}$ the num-
36 ber of subjects classified into category i by rater 1 and $n_{\cdot j}$ be the number of
37 subjects classified into category j by rater 2. Denote by p_{ij} , $p_{i\cdot}$ and $p_{\cdot j}$ the
38 corresponding proportions ($i, j = 1, \dots, K$).

39

Table 1

Two-way contingency table resulting from the classification of n items by 2 raters on an ordinal scale with K categories

		Rater 2					
Rater 1	1	...	j	...	K	Total	
1	n_{11}	...	n_{1j}	...	n_{1K}	$n_{1.}$	
i	n_{i1}	...	n_{ij}	...	n_{iK}	$n_{i.}$	
K	n_{K1}	...	n_{Kj}	...	n_{KK}	$n_{K.}$	
Total	$n_{.1}$...	$n_{.j}$...	$n_{.K}$	n	

The weighted kappa coefficient can be defined in terms of agreement weights by

$$\kappa_w = \frac{p_o - p_e}{1 - p_e} \quad (1)$$

where $p_o = \sum_{i=1}^K \sum_{j=1}^K w_{ij} p_{ij}$ and $p_e = \sum_{i=1}^K \sum_{j=1}^K w_{ij} p_{i.} p_{.j}$ with $0 \leq w_{ij} \leq 1$ and $w_{jj} = 1$ ($i, j = 1, \dots, K$), or in terms of disagreement weights by

$$\kappa_w = 1 - \frac{q_o}{q_e} \quad (2)$$

where $q_o = \sum_{i=1}^K \sum_{j=1}^K v_{ij} p_{ij}$ and $q_e = \sum_{i=1}^K \sum_{j=1}^K v_{ij} p_{i.} p_{.j}$ with $0 \leq v_{ij} \leq 1$ and $v_{jj} = 0$ ($i, j = 1, \dots, K$). However, the weighted kappa coefficient can also be obtained using unscaled disagreement weights, i.e., v_{ij} not restricted to the $[0,1]$ interval.

44

45 Cohen's kappa coefficient is a particular case of the weighted kappa coefficient where $w_{ij} = 1$ ($v_{ij} = 0$) for $i = j$ and $w_{ij} = 0$ ($v_{ij} = 1$) for $i \neq j$

47 $(i, j = 1, \dots, K)$. Cicchetti and Allison (1971) proposed "linear" weights of
48 the form $w_{ij} = 1 - |i - j|/(K - 1)$, whereas Fleiss and Cohen (1973) used
49 the quadratic weights $w_{ij} = 1 - (i - j)^2/(K - 1)^2$. The disagreement weights
50 $v_{ij} = (i - j)^2$ are also commonly used (Ludbrook (2002); Agresti (2002)) as
51 are the linear disagreement weights $v_{ij} = |i - j|$.

52

53 Cohen (1968) showed that if the marginal distributions of the 2 raters are
54 the same and if the weights of disagreement are defined as $v_{ij} = (i - j)^2$, the
55 weighted kappa coefficient is equivalent to the product-moment correlation co-
56 efficient. Furthermore, Fleiss and Cohen (1973) showed that using the weights
57 v_{ij} , the weighted kappa coefficient has the same interpretation as the intra-
58 class correlation coefficient of reliability when systematic variability between
59 raters is included as a component of total variation. More recently, Schus-
60 ter (2004) explicitly decomposed the weighted kappa coefficient defined with
61 the quadratic disagreement weights in terms of rater means, rater variances
62 and rater covariance in the context of a two-way analysis of variance. To the
63 best of our knowledge, no interpretation was given for the weighted agreement
64 coefficient with linear agreement or disagreement weights.

65 **3 THE REVISITED WEIGHTED KAPPA COEFFICIENT**

Hereafter, we shall focus on the linear weights introduced by Cicchetti and Allison (1971) ($w_{ij} = 1 - |i - j|/(K - 1)$) and revisit the weighted kappa coefficient for an ordinal scale. The interpretation of the agreement index obtained with the linear disagreement weights ($v_{ij} = |i - j|$) will follow straightforwardly

Table 2

Reduction of the $K \times K$ contingency table into a 2×2 classification table by selecting a cut-off level k ($k = 1, \dots, K$) on the ordinal scale (see text)

Rater 2			
Rater 1	$\leq k$	$> k$	Total
$\leq k$	$N_{11}(k)$	$N_{12}(k)$	$N_{1.}(k)$
$> k$	$N_{21}(k)$	$N_{22}(k)$	$N_{2.}(k)$
Total	$N_{.1}(k)$	$N_{.2}(k)$	n

since

$$w_{ij} = 1 - \frac{v_{ij}}{K - 1}. \quad (3)$$

66 For any "cut-off" value k ($k = 1, \dots, K - 1$), the $K \times K$ contingency table
 67 (see Table 1) can be reduced into a 2×2 classification table by summing up all
 68 observations below and above the first k rows and first k columns (see Table
 69 2) where

$$N_{11}(k) = \sum_{i=1}^k \sum_{j=1}^k n_{ij} \quad N_{12}(k) = \sum_{i=1}^k \sum_{j=k+1}^K n_{ij}$$

$$N_{21}(k) = \sum_{i=k+1}^K \sum_{j=1}^k n_{ij} \quad N_{22}(k) = \sum_{i=k+1}^K \sum_{j=k+1}^K n_{ij}$$

Let $F_{lm}(k) = \frac{1}{n} N_{lm}(k)$, $F_{l.} = \frac{1}{n} N_{l.}(k)$ and $F_{.m} = \frac{1}{n} N_{.m}(k)$ be the corresponding joint and marginal frequencies ($l, m = 1, 2; k = 1, \dots, K - 1$). Finally, denote by

$$p_o(k) = F_{11}(k) + F_{22}(k) \quad (4)$$

and

$$p_e(k) = F_{1.}(k)F_{.1}(k) + F_{2.}(k)F_{.2}(k) \quad (5)$$

70 the observed and expected weighted agreements corresponding to Table 2.

Now, consider the quantities

$$p_o^* = \frac{1}{K-1} \sum_{k=1}^{K-1} p_o(k) \quad (6)$$

and

$$p_e^* = \frac{1}{K-1} \sum_{k=1}^{K-1} p_e(k) \quad (7)$$

71 We show that $p_o^* = p_o$ and $p_e^* = p_e$ where p_o and p_e are respectively the "lin-
72 early" weighted proportions of observed and expected agreement, as defined
73 by Cicchetti & Allison (1971).

74

75 Since

76

$$\begin{aligned} p_o^* &= \frac{1}{K-1} \sum_{k=1}^{K-1} \left(\sum_{i=1}^k \sum_{j=1}^k p_{ij} + \sum_{i=k+1}^K \sum_{j=k+1}^K p_{ij} \right) \\ &= \frac{1}{K-1} \sum_{k=1}^{K-1} \left(\sum_{i=1}^K \sum_{j=1}^K p_{ij} - \sum_{i=1}^k \sum_{j=k+1}^K p_{ij} - \sum_{i=k+1}^K \sum_{j=1}^k p_{ij} \right) \\ &= \sum_{i=1}^K \sum_{j=1}^K p_{ij} - \frac{1}{K-1} \sum_{k=1}^{K-1} \left(\sum_{i=1}^k \sum_{j=k+1}^K p_{ij} + \sum_{i=k+1}^K \sum_{j=1}^k p_{ij} \right) \end{aligned} \quad (8)$$

77 and

78

$$\begin{aligned} p_o &= \sum_{i=1}^K \sum_{j=1}^K \left(1 - \frac{|i-j|}{K-1} \right) p_{ij} \\ &= \sum_{i=1}^K \sum_{j=1}^K p_{ij} - \frac{1}{K-1} \sum_{i=1}^K \sum_{j=1}^K |i-j| p_{ij} \\ &= \sum_{i=1}^K \sum_{j=1}^K p_{ij} - \frac{1}{K-1} \sum_{i=1}^K \sum_{j=1}^i (i-j) p_{ij} - \frac{1}{K-1} \sum_{i=1}^{K-1} \sum_{j=i+1}^K (j-i) p_{ij}, \end{aligned} \quad (9)$$

it suffices to prove that

$$\sum_{k=1}^{K-1} \left(\sum_{i=1}^k \sum_{j=k+1}^K p_{ij} + \sum_{i=k+1}^K \sum_{j=1}^k p_{ij} \right) = \sum_{i=1}^{K-1} \sum_{j=i+1}^K (j-i)p_{ij} + \sum_{i=1}^K \sum_{j=1}^i (i-j)p_{ij} \quad (10)$$

79 We have successively,

80

$$\begin{aligned} & \sum_{k=1}^{K-1} \left(\sum_{i=1}^k \sum_{j=k+1}^K p_{ij} + \sum_{i=k+1}^K \sum_{j=1}^k p_{ij} \right) = \sum_{k=1}^{K-1} \sum_{i=1}^k \sum_{j=k+1}^K p_{ij} + \sum_{k=1}^{K-1} \sum_{i=k+1}^K \sum_{j=1}^k p_{ij} \\ &= \sum_{i=1}^1 \sum_{j=2}^K p_{ij} + \sum_{i=1}^2 \sum_{j=3}^K p_{ij} + \cdots + \sum_{i=1}^{K-1} \sum_{j=K}^K p_{ij} \\ &+ \sum_{i=2}^K \sum_{j=1}^1 p_{ij} + \sum_{i=3}^K \sum_{j=1}^2 p_{ij} + \cdots + \sum_{i=K}^K \sum_{j=1}^{K-1} p_{ij} \\ &= \sum_{i=1}^1 \sum_{j=2}^K p_{ij} + \sum_{i=1}^1 \sum_{j=3}^K p_{ij} + \sum_{i=2}^2 \sum_{j=3}^K p_{ij} + \cdots + \sum_{i=1}^1 \sum_{j=K}^K p_{ij} + \sum_{i=2}^{K-1} \sum_{j=K}^K p_{ij} \\ &+ \sum_{i=K}^K \sum_{j=1}^1 p_{ij} + \sum_{i=2}^{K-1} \sum_{j=1}^1 p_{ij} + \sum_{i=K}^K \sum_{j=1}^2 p_{ij} + \sum_{i=3}^{K-1} \sum_{j=1}^2 p_{ij} + \cdots + \sum_{i=K}^K \sum_{j=1}^{K-1} p_{ij} \\ &= \sum_{j=2}^K (j-1)p_{1j} + \sum_{i=2}^2 \sum_{j=3}^K p_{ij} + \cdots + \sum_{i=2}^{K-1} \sum_{j=K}^K p_{ij} \\ &+ \sum_{j=1}^{K-1} (K-j)p_{Kj} + \sum_{i=2}^{K-1} \sum_{j=1}^1 p_{ij} + \cdots + \sum_{i=3}^{K-1} \sum_{j=1}^2 p_{ij} \\ &= \sum_{j=2}^K (j-1)p_{1j} + \sum_{j=3}^K (j-2)p_{2j} + \cdots + \sum_{j=K}^K (j-(K-1))p_{K-1,j} \\ &+ \sum_{j=1}^{K-1} (K-j)p_{Kj} + \sum_{j=1}^{K-2} (K-1-j)p_{K-1,j} + \cdots + \sum_{j=1}^{K-(K-1)} (K-(K-1)-j)p_{K-(K-1),j} \\ &= \sum_{i=1}^{K-1} \sum_{j=i+1}^K (j-i)p_{ij} + \sum_{i=1}^K \sum_{j=1}^i (i-j)p_{ij} \end{aligned} \quad (11)$$

81 Thus, $p_o^* = p_o$. The proof for $p_e^* = p_e$ proceeds similarly. Thus, using the
82 linear agreement weights introduced by Cicchetti and Allison (1971), the ob-
83 served and expected weighted agreements are merely the mean values of the
84 corresponding proportions of all 2×2 tables obtained by collapsing the first
85 k categories and last $K - k$ categories ($k = 1, \dots, K - 1$) of the original
86 $K \times K$ classification table. When considering the linear disagreement weights,

87 the observed and expected weighted disagreements correspond to the sum of
88 the observed and expected proportions of disagreement of the $K - 1$ embedded
89 2×2 tables, respectively.

90 4 EXAMPLE

91 Gilmour et al. (1997) conducted an agreement study to compare two meth-
92 ods for assessing cervical ectopy, defined as the presence of endocervical-type
93 columnar epithelium on the portio surface of the cervix. A computerized
94 planimetry method was developed for measuring cervical ectopy, and the re-
95 liability of that method was compared with direct visual assessment. Pho-
96 tographs of the cervix of 85 women without cervical disease were assessed for
97 cervical ectopy by three medical raters who used both assessment methods.
98 The response of interest, cervical ectopy size, was an ordinal variable with
99 four categories: (1) minimal, (2) moderate, (3) large and (4) excessive. The
100 contingency table for two of the three raters using the visual method is dis-
101 played in Table 3. In each cell, the first term corresponds to the cell count,
102 the second term to the linear agreement weight and the third one to the linear
103 disagreement weight.

104

105 When computing the weighted observed and expected agreements, we obtain
106 $p_o = 0.800$, $p_e = 0.583$, yielding $\kappa_w = 0.520$. Since $K = 4$, three "embedded"
107 2×2 tables can be constructed as described before (see Table 4). From these
108 tables, we calculate $p_o^* = \frac{1}{3} \sum_{k=1}^3 p_o(k) = (0.812 + 0.788 + 0.800)/3 = 0.800$
109 and $p_e^* = \frac{1}{3} \sum_{k=1}^3 p_e(k) = (0.618 + 0.506 + 0.626)/3 = 0.583$. These are as
110 expected equal to p_o and p_e , respectively. It should be remarked that the av-

Table 3

Two-way contingency table resulting from cervical ectopy ratings using the visual method by two raters

Rater 1	Rater 2				Total
	1	2	3	4	
1	13 ^a	2	0	0	15
	1.0 ^b	0.67	0.33	0.0	
	0.0 ^c	1.0	2.0	3.0	
2	10	16	3	0	29
	0.67	1.0	0.67	0.33	
	1.0	0.0	1.0	2.0	
3	3	7	3	0	13
	0.33	0.67	1.0	0.67	
	2.0	1.0	0.0	1.0	
4	1	4	12	11	28
	0.0	0.33	0.67	1.0	
	3.0	2.0	1.0	0.0	
Total	27	29	18	11	85

^a Observed counts

^b Linear agreement weights $w_{ij} = 1 - |i - j| / (K - 1)$

^c Linear disagreement weights $v_{ij} = |i - j|$

111 erage kappa coefficient derived from the tables, namely $\bar{\kappa} = \frac{1}{3} \sum_{k=1}^3 \kappa(k) =$
112 $(0.507 + 0.572 + 0.465)/3 = 0.515$, differs from κ_w . The weighted observed
113 and expected disagreements are equal to $q_o = 0.600$ and $q_e = 1.25$, respec-
114 tively, yielding a weighted kappa coefficient of $\kappa_w = 0.52$. From the embed-
115 ded tables, we have $q_o^* = \sum_{k=1}^3 q_o(k) = 0.188 + 0.212 + 0.200 = 0.600$ and
116 $p_e^* = \sum_{k=1}^3 q_e(k) = 0.382 + 0.494 + 0.374 = 1.25$, as expected.

117

118 5 DISCUSSION

119 The weighted kappa coefficient is widely used to quantify the agreement be-
120 tween 2 raters on an ordinal scale. The weights are generally given a priori and
121 defined arbitrarily. Graham and Jackson (1993) observed that the value of the
122 weighted kappa coefficient can vary considerably according to the weighting
123 scheme used and henceforth may lead to different conclusions. In practice, the
124 linear (Cicchetti and Allison, 1971) and quadratic (Fleiss and Cohen, 1973)
125 weighting schemes are the most widely used. Quadratic weights have received
126 much attention in the literature because of their practical interpretation. For
127 instance, Fleiss and Cohen (1973) and Schuster (2004) showed that using the
128 weights $v_{ij} = (i - j)^2$, the weighted kappa coefficient can be interpreted as
129 an intraclass correlation coefficient in a two-way analysis of variance setting.
130 In this article, we focused on the linearly weighted kappa coefficient defined
131 by Cicchetti and Allison (1971) or equivalently defined by the linear disagree-
132 ment weights $v_{ij} = |i - j|$ and strove to give an intuitive interpretation of it.
133 Specifically, we showed that the observed and expected weighted agreements
134 are merely the mean values of the corresponding proportions of all 2×2

Table 4

All possible embedded 2×2 classification tables ($k = 1, 2, 3$) derived from the original 4×4 contingency table for cervical ectopy ratings by two raters

Rater 2				Rater 2			
Rater 1	≤ 1	> 1	Total	Rater 1	≤ 2	> 2	Total
≤ 1	13	2	15	≤ 2	41	3	44
> 1	14	56	70	> 2	15	26	41
Total	27	58	85	Total	56	29	85

$$p_o(1) = 0.812; q_o(1) = 0.188$$

$$p_o(2) = 0.788; q_o(2) = 0.212$$

$$p_e(1) = 0.618; q_e(1) = 0.382$$

$$p_e(2) = 0.506; q_e(2) = 0.494$$

$$\kappa(1) = 0.507$$

$$\kappa(2) = 0.572$$

Rater 2			
Rater 1	≤ 3	> 3	Total
≤ 3	57	0	57
> 3	17	11	28
Total	74	11	85

$$p_o(3) = 0.800; q_o(3) = 0.200$$

$$p_e(3) = 0.626; q_e(3) = 0.374$$

$$\kappa(3) = 0.465$$

135 tables obtained by collapsing the first k categories and last $K - k$ categories
136 ($k = 1, \dots, K - 1$) of the original $K \times K$ classification table. It should be noted,
137 however, that the weighted agreement coefficient derived from the original ta-
138 ble is not equal to the mean value of the non-weighted $K - 1$ κ coefficients
139 obtained from the 2×2 collapsed tables. When using linear disagreement
140 weights, the weighted observed and expected disagreements are obtained by
141 the sum rather than the average of the corresponding elements of the 2×2
142 tables. In other words, the linearly weighted kappa coefficient can simply be
143 derived from $K - 1$ embedded 2×2 classification tables. The linear form of
144 the kappa coefficient, besides its simplicity, presents some advantages over the
145 quadratic version. As demonstrated by Brenner and Kliebach (1996), it is less
146 sensitive to the number of categories and should therefore be preferred when
147 the number of categories of the ordinal scale is large. As a conclusion, we have
148 shown that the linearly weighted kappa coefficient for a K -ordinal table can
149 be naturally derived from non-weighted observed and expected agreements
150 (disagreements) computed from $K - 1$ embedded 2×2 classification tables.

151 **References**

- [1]₃₂ Agresti, A., 2002. Categorical data analysis (2nd ed.). New York: John Wiley
153 and Sons.
- [2]₄₄ Brenner, H., Kliebach, U., 1996. Dependence of weighed kappa coefficients on
155 the number of categories. *Epidemiology* 7, 199–202.
- [3]₃₆ Cicchetti, D., Allison, T., 1971. A new procedure for assessing reliability of
157 scoring eeg sleep recordings. *American Journal EEG Technology* 11, 101–
158 109.
- [4]₃₉ Cohen, J., 1960. A coefficient of agreement for nominal scales. *Educational*

160 and Psychological Measurement 20, 37–46.

[5]₁ Cohen, J., 1968. Weighted kappa: nominal scale agreement with provision for
162 scaled disagreement or partial credit. Psychological Bulletin 70, 213–220.

[6]₃ Fleiss, J. L., Cohen, J., 1973. The equivalence of weighted kappa and the
164 intraclass correlation coefficient as measure of reliability. Educational and
165 Psychological Measurement 33, 613–619.

[7]₆ Gilmour, E., Ellerbrock, T., Koulos, J., Chiasson, M., Williamson, J., Kubn,
167 L., Wright, T. J., 1997. Measuring cervical ectopy: direct visual assessment
168 versus computerized planimetry. American Journal of Obstetrics and Gy-
169 necology 176, 108–111.

[8]₀ Graham, P., Jackson, R., 1993. The analysis of ordinal agreement data: beyond
171 weighted kappa. Journal of Clinical Epidemiology 46, 1055–1062.

[9]₂ Ludbrook, J., 2002. Statistical techniques for comparing measurers and meth-
173 ods of measurement: a critical review. Clinical and Experimental Pharma-
174 cology and Physiology 29, 527–536.

[10] Schuster, C., 2004. A note on the interpretation of weighted kappa and its
176 relation to other rater agreement statistics for metric scales. Educational
177 and Psychological Measurement 64, 243–253.