

## Measurement of Rater Consistency by Chance-Corrected Agreement Coefficients

Zheng Xie

School of Engineering  
University of Central Lancashire  
Preston, UK  
e-mail: zxie2@uclan.ac.uk

Chaitanya Gadepalli

University Department of Otolaryngology  
Central Manchester University Hospitals Foundation  
Trust and University of Manchester Academic  
Health Science Centre  
Manchester, UK  
e-mail: cgadepalli@gmail.com

Barry M.G. Cheetham

School of Computer Science  
University of Manchester  
Manchester, UK  
e-mail: barry@manchester.ac.uk

**Abstract**— Measurement of consistency in the decisions made by observers or raters is an important problem in clinical medicine. Chance corrected agreement coefficients such as the Cohen and Fleiss Kappas are commonly used for this purpose, though the way that they estimate the probability of agreement 'by chance' has been strongly questioned. Alternatives have been proposed, such as the Aickin Alpha coefficient and the Gwet AC<sub>1</sub> and AC<sub>2</sub> coefficients which are gaining currency. A well known paradox illustrates deficiencies of the Kappa coefficients which, it is claimed, are remedied by an approach which grades the subjects according to their probability of being hard to score. The AC<sub>1</sub> and AC<sub>2</sub> coefficients result from the application of this grading to the Brennan-Prediger coefficient which may be considered a simplified form of Kappa. This paper questions the rationale of the hardness probability used by AC<sub>1</sub> and proposes an alternative approach that may be applied to weighted and unweighted multi-rater Cohen and Fleiss Kappas and also Intra-Class Correlation (ICC) coefficients.

**Keywords**- Rater consistency, Fleiss Kappa, Cohen Kappa, ICC, Gwet's AC<sub>1</sub> coefficient

### I. INTRODUCTION

The consistency of clinical observations is clearly important in the practice of medicine. In such practice, a patient may see just one clinician. However, it may be useful to have the reassurance that the observations are likely to be independent of the choice of clinician since these observations are likely to affect diagnoses and ultimate treatment. If the observations are found to be different for different clinicians, this may not necessarily be a bad thing as different perspectives can be valuable. However, it would be useful to know to what extent assessments are likely to be consistent and whether certain phenomena are

hard or easy to assess. Such knowledge may, for example, suggest when a second opinion may be valuable.

Investigating how much consistency is likely, and finding ways of improving it requires a clinical trial with a selection of subjects comprising patients and other volunteers and a number of clinical observers referred to as raters. Such a trial was carried out for a voice quality assessment procedure with 102 subjects and five raters [1]. This trial required measurements of the intra-rater (self) consistency of decisions by the same raters at different times and also of inter-rater consistency of decisions by different raters observing the same subjects.

The decisions may be diagnoses of medical conditions or the severity of such conditions. The decisions may be categorical and denoted by labels. Or they may be ordinal which means that they are numbers often referred to as scores. Unlike labels, scores have magnitudes and may be compared in terms of differences between them.

Given  $N$  subjects and  $R$  raters who each observe all subjects, the 'proportion of agreement'  $P_o$  may be considered as a measure of consistency for categorical or ordinal decisions [2]. Denoting by  $A(i,r)$  the decision or score given by rater  $r$  to subject  $i$ ,  $P_o$  may be expressed as:

$$P_o = \frac{1}{NL} \sum_{j=1}^N \sum_{r=1}^R \sum_{s=r+1}^R \delta(A(i,r), A(i,s)) \text{ where } \delta(u,v) = \begin{cases} 1: u = v \\ 0: u \neq v \end{cases} \quad (1)$$

In equation (1),  $L = R(R-1)/2$  which is the number ways in which two distinct raters,  $r$  and  $s$ , may be selected from the  $R$  raters for comparison of their decisions. Therefore  $P_o$  counts the number of times that a pair of raters agree for a subject. It is expressed as a proportion of the number of subjects times the number of rater pairs, and is a number between 0 and 1. When there is complete agreement by all raters for all subjects,  $P_o$  will be equal to 1. When there is almost no agreement,  $P_o$  will be close to zero.

As defined for categorical (nominal) decisions,  $P_o$  may be used also for ordinal scoring if the scores are considered as labels rather than numbers. However this gives equal weight to all possible differences in scores regardless of their magnitudes. For ordinal scoring it is often preferable to give more importance to larger differences than small difference, and this leads to a weighted version of  $P_o$  defined as follows:

$$P_o = \frac{1}{NL} \sum_{i=1}^N \sum_{r=1}^R \sum_{s=r+1}^R w(A(i,r), A(i,s)) \quad (2)$$

where  $w(u,v)$  is a weighting function [3]. Assuming that there are  $Q$  possible scores  $1, 2, \dots, Q$ , for linear weighting:

$$w(u,v) = 1 - C(u,v) \quad \text{where } C(u,v) = |u - v| / (Q - 1) \quad (3)$$

for quadratic weighting,

$$w(u,v) = 1 - C(u,v) \quad \text{where } C(u,v) = (u - v)^2 / (Q - 1)^2 \quad (4)$$

and for no weighting,

$$w(u,v) = 1 - C(u,v) \quad \text{where } C(u,v) = 1 - \delta(u,v) \quad (5)$$

When  $w(u,v)$  is defined by equation (5), equation (2) becomes identical to equation (1) as applied to ordinal scores considered as labels. There are many other possible weighting-functions that may be considered, but these three are of special interest. In equations (3-5),  $C(u,v)$  is a cost function which determines the degree to which the value of  $P_o$  is decreased from unity by a rater-pair disagreement. With linear weighting, the magnitude of the score difference, scaled by  $(Q-1)$ , constitutes the decrease or cost of the disagreement. With quadratic weighting, these magnitudes are squared and scaled by  $(Q-1)^2$  to produce the cost. With no weighting, any score difference contributes the same unit cost. As with the unweighted version of  $P_o$ , the weighted version is equal to 1 for perfect agreement. The scaling of the cost by  $(Q-1)$  for linear weighting and  $(Q-1)^2$  for quadratic weighting makes the minimum possible value of  $P_o$  equal to 0 which would occur when all rater-pairs disagree to the maximum possible extent for all subjects.

## II. CHANCE-CORRECTED AGREEMENT COEFFICIENTS

Unweighted and weighted versions of  $P_o$  are straightforward measures of consistency for categorical or ordinal scoring. But they are biased by the probability of some agreement occurring by chance. If all raters were to make random decisions evenly distributed over  $Q$  categories or scores, 'by chance' agreement would be expected with a probability of  $1/Q$ , even if the raters made their decisions without even seeing the subjects. This would make unweighted  $P_o$  equal to  $1/Q$  and weighted  $P_o$  equal to  $T_w/Q^2$  [3] where:

$$T_w = \sum_{k=1}^Q \sum_{\ell=1}^Q w(k,\ell) \quad (6)$$

With four scoring categories and 'by chance' scoring, the expectation of unweighted  $P_o$  would be  $1/4$  or 25% for an even spread of decisions over the four categories, and with an uneven spread of decisions,  $P_o$  could be even greater, thus giving a false impression of some consistency when there may be none.

Chance corrected agreement coefficients aim to cancel out the bias in  $P_o$ , while still providing a number between 0 and 1. They are normally expressed as:

$$\gamma = \frac{P_o - P_e}{1 - P_e} \quad (7)$$

where  $P_o$  is as defined above and  $P_e$  is an estimate of the probability of agreement by chance. If there is almost complete agreement,  $P_o$  will be close to 1 and  $\gamma$  will be close to 1 unless  $P_e$  is also close to 1. If  $P_e$  is close to 1, almost all agreement would be considered to have occurred by chance.

## III. BRENNAN-PREDIGER COEFFICIENT

The simplest chance corrected agreement coefficient is known as the Brennan-Prediger coefficient [4]. The unweighted or categorical version is:

$$\gamma = \frac{P_o - P_e}{1 - P_e} \quad \text{where } P_e = \frac{1}{Q} \quad (8)$$

with  $P_o$  defined by equation (1). The weighted version for ordinal scoring is:

$$\gamma = \frac{P_o - P_e}{1 - P_e} \quad \text{where } P_e = \frac{T_w}{Q^2} \quad (9)$$

with  $P_o$  defined weighted by equation (2) and  $T_w$  by equation (6). The Brennan-Prediger coefficient cancels out the bias in  $P_o$  for an even distribution of rater decisions or scores among the  $Q$  categories. But it will not do this accurately for uneven distributions.

## IV. MULTI-RATER COHEN KAPPA

The Cohen Kappa aims to remove the bias present in  $P_o$  by estimating and taking into account the probability of agreement 'by chance' given the distribution of decisions produced by each rater. It was originally proposed [5] for categorical rating by two raters, but was generalised by Hubert [6] and Conger [7] to a multi-rater version. In categorical form, it may be expressed [2] as follows:

$$\gamma = \frac{P_o - P_e}{1 - P_e} \quad \text{with } P_e = \frac{1}{LN^2} \sum_{i=1}^N \sum_{j=1}^N \sum_{r=1}^R \sum_{s=r+1}^R \delta(A(i,r), A(j,s)) \quad (10)$$

and with  $P_o$  defined by equation (1). In equation (10),  $P_e$  is an estimate of the probability of a pair of arbitrary raters agreeing by chance when these raters make arbitrary decisions which are not the same as those of the actual raters, but are similarly distributed over the  $Q$  possible categories or scores.

For  $P_e$  to be a reliable estimate of 'by chance' scoring for a population of subjects, it has to be assumed that the  $N$  subjects are a reasonable sample of the population. If there is a bias towards one particular score in the scores obtained for the  $N$  subjects, then it is assumed that that the same bias exists in the population. Otherwise the sample will be unrepresentative of the population and the estimate of  $P_e$  will be unreliable.

Equation (10) becomes identical to the original Cohen Kappa when the number of raters,  $R$ , is equal to two. The generalisation to more than two raters is due to Conger [7] and Hubert [6]. An alternative generalisation by Light [8] is also identical for two raters but slightly different for more. Equation (10) may be further generalised to weighted form [9] for ordinal scoring by defining  $P_o$  by equation (2) and  $P_e$  by equation (11):

$$P_e = \frac{1}{LN^2} \sum_{i=1}^N \sum_{j=1}^N \sum_{r=1}^R \sum_{s=r+1}^R w(A(i,r), A(j,s)) \quad (11)$$

As with the Brennan-Prediger coefficient, the introduction of weighting increases the cost of any by chance disagreement between arbitrary raters. But  $P_e$  is no longer a constant independent of the scoring distribution. Combining equations (2), (7) and (11) we obtain an expression for the weighted multi-rater Cohen Kappa which can be re-expressed [2] as follows:

$$\gamma = 1 - \frac{(1/N) \sum_{r=1}^R \sum_{s=r+1}^R \sum_{i=1}^N C(A(i,r), A(i,s))}{(1/N^2) \sum_{r=1}^R \sum_{s=r+1}^R \sum_{i=1}^N \sum_{j=1}^N C(A(i,r), A(j,s))} \quad (12)$$

## V. FLEISS-KAPPA

The Fleiss Kappa [10] is an alternative chance-corrected agreement coefficient originally defined for two or more categorical raters. When used for two raters, it becomes identical to the Scott Pi coefficient [11]. A Fleiss Kappa of 1 indicates perfect agreement between all raters, and lower values are interpreted on a scale similar to that assumed for the Cohen Kappa.

For categorical scoring with  $R$  raters and  $Q$  scoring categories, Fleiss [10] calculates  $P_o$ , by equation (1), and then calculates the proportion,  $\pi_k$ , of all assignments to category  $k$ , for all raters and all subjects, for  $k = 1, 2, \dots, Q$ , as follows:

$$\pi_k = \frac{1}{NR} \sum_{i=1}^N \sum_{r=1}^R \delta(A(i,r), k) \quad (13)$$

If it is assumed that the  $N$  subjects are a reasonable sample of the population,  $\pi_k$  may be considered an estimate of the probability that a randomly selected rater will classify a randomly selected subject into category  $k$ . As there are  $Q$  categories, it follows that:

$$\sum_{k=1}^Q \pi_k = 1 \quad (14)$$

Fleiss [10] then estimates the probability of agreement 'by chance' as:

$$P_e = \sum_{k=1}^Q \pi_k^2 \quad (15)$$

A weighted version [3] of the Fleiss Kappa is obtained by defining  $P_o$  by equation (2) and  $P_e$  by equation (16):

$$P_e = \sum_{k=1}^Q \sum_{\ell=1}^Q w(k, \ell) \pi_k \pi_\ell \quad (16)$$

Substituting these values of  $P_o$  and  $P_e$  into the Kappa equation (3) gives an expression for Fleiss Kappa which can be re-expressed [2] as:

$$\gamma = 1 - \frac{\frac{1}{NL} \sum_{r=1}^R \sum_{s=r+1}^R \sum_{i=1}^N C(A(i,r), A(i,s))}{\frac{1}{(Nn)^2} \sum_{r=1}^R \sum_{s=1}^R \sum_{i=1}^N \sum_{j=1}^N C(A(i,r), A(j,s))} \quad (17)$$

where the cost function  $C$  may apply no weighting, or linear, quadratic or other weighting as in equations (3-5). This equation is valid for any number,  $R$ , of raters including  $R = 2$ .

The Fleiss and Cohen versions of Kappa differ because Fleiss specifies that each rater index does not necessarily refer to the same person. In this case, it is considered inappropriate to define  $P_e$  in terms of the distribution of scores for each rater index. The more general assumption made by Fleiss about the likely distribution of scores for each rater index, i.e. that they are all equal, seems more appropriate. The Fleiss Kappa [10] is unaffected by the characteristic trends in the scoring by individuals. Only the distribution of scores among the  $Q$  scoring categories by all raters is considered important. The differences between the Fleiss Kappa and the multi-rater Cohen Kappa are often small but sometimes noticeable. Where there are individual (fixed) raters it is probably best to use the multi-rater Cohen Kappa rather than the Fleiss Kappa. This is because it preserves the original definition of agreement by chance, which takes into account the typical scoring distributions of the individual raters.

## VI. MISSING SCORES

The equations given above for the multi-rater Cohen, Fleiss and Brennan-Prediger coefficients assume that all  $N$  subjects are scored by all  $R$  raters. They have generalised by Gwet [3] to the case where some scores are missing, but we do not consider this case here.

## VII. GWET'S PARADOX

There is controversy about the way Cohen [5] and Fleiss Kappa [10] estimate by chance agreement (14) and different approaches, such as the  $AC_1$  and  $AC_2$  coefficients by Gwet [3], are gaining currency. The deficiencies of the Cohen and Fleiss Kappas are illustrated by the example in Table 1, which is similar to examples quoted by Gwet [3].

In this example, there are two raters for 20 subjects with two possible scoring categories. Rater 1 scores all subjects in category 1, and rater 2 scores 18 out of 20 in category 1. It may appear that that there is a high level of agreement. However, since  $P_o = 0.9$  and  $P_e = 0.9$  for unweighted Cohen Kappa and  $P_e = 0.905$  for unweighted Fleiss Kappa, both Kappas give zero or a value close to zero, thus indicating little or no agreement. The problem lies with the estimation of  $P_e$ , since almost all agreement is classified as agreement by chance. This happens because the 20 subjects are not representative of a population considered typical of subjects in general. By equation (8), the Brennan-Prediger coefficient, which is independent of the distribution of scores, gives the more reasonable value of 0.8.

Ideally when defining a chance-corrected agreement measure we should specify the expected scoring characteristics for a population of subjects. If this is a population for which the overwhelming majority of subjects are expected to be scored as category 1, then the Cohen and Fleiss Kappas defined above may be adequate. However we normally do not expect such a population, and are more likely to expect a more even distribution of scores. By default, an assumed distribution of population scores could be an even distribution among the  $Q$  categories. For such a distribution, the Brennan-Prediger coefficient is appropriate, and neither the Cohen nor the Fleiss Kappa are appropriate when the actual scores are not evenly distributed. However, we may have reason to expect some scores to occur more frequently than others, and the evidence may lie in the scores themselves. This is where the Cohen and Fleiss Kappas may prove useful, though some modifications are needed.

## VIII. FURTHER INVESTIGATION OF GWET'S PARADOX

The paradox illustrated above, and referred to by Gwet [3], occurs whenever the value of  $\pi_k$  becomes close to 1 for some value of  $k$ . In this case, by equation (14), all other values of  $\pi_k$  become close to zero. To investigate this situation, we randomly generated a set of scores for  $N = 50$  subjects,  $R = 5$  raters and  $Q = 4$  scoring categories. Given  $Q$  values of  $\pi_k$ , we generated a random score in the range 1 to  $Q$  (inclusive) for each subject index  $i$ , for each of the  $R$

raters, such that the overall probability of getting score  $k$  was equal to  $\pi_k$  for  $k = 1, 2, \dots, Q$ .

Initially, we made  $\pi_k = 1/Q$  for  $k = 1, 2, \dots, Q$ , which meant that all scores were equally probable over all subjects and all raters. By equations (10 and (15), this case gives  $P_e = 1/Q$  for both the unweighted Cohen and Fleiss Kappas, which are lowest possible values.  $P_e$  is always equal to  $1/Q$  for the unweighted Brennan-Prediger coefficient.

We then randomly generated further sets of random scores, with one of the  $\pi_k$  values increased from  $1/Q$ , and the other  $(Q-1)$  values decreased to satisfy equation (15). We chose  $\pi_1$  to be the value that increased, and made all other  $\pi_k$  values equal to  $(1 - \pi_1)/(Q-1)$ . By gradually increasing  $\pi_1$  towards 1 we generated a series of scoring patterns that gradually approached the maximally concentrated distribution where  $\pi_1 = 1$  and all other values of  $\pi_k$  are zero. For this maximally concentrated distribution, all raters give score 1, to all subjects.

The situation when  $\pi_1$  becomes close to 1 further demonstrates the paradox pointed out by Gwet [3]. The resulting values of unweighted Cohen and Fleiss Kappas are plotted against increasing  $\pi_1$  in Figure 1, along with the Brennan-Prediger coefficient and the Gwet  $AC_1$  coefficient. It may be seen in Figure 1 that as  $\pi_1$  approaches 1, both the Cohen and Fleiss Kappas remain close to zero indicating no agreement except by chance. The Brennan Prediger coefficient approaches 1 (perfect agreement) as  $\pi_1$  approaches 1.

The corresponding values of  $P_e$  for each of the coefficients plotted in Figure 1 are plotted against  $\pi_1$  in Figure 2. Perhaps unexpectedly, the probability of agreement by chance, as estimated by both the unweighted Cohen and Fleiss Kappas, increases as the scores become more and more concentrated on score 1.  $P_e$  for the unweighted Brennan-Prediger coefficient remains constant at  $1/Q$  with  $Q=4$ . The Gwet  $AC_1$  coefficient will be discussed in the next Section.

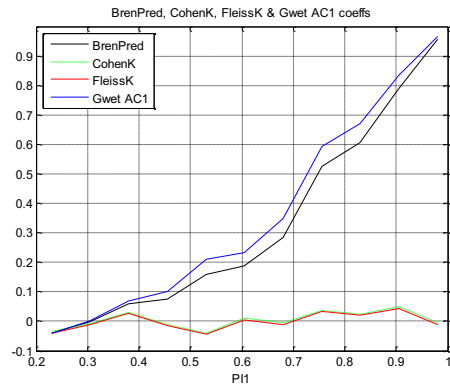


Figure 1. Graph comparing unweighted Cohen and Fleiss Kappas with the Brennan-Prediger and Gwet's  $AC_1$  coefficients for increasing concentration of scores (Kappas almost coincide).

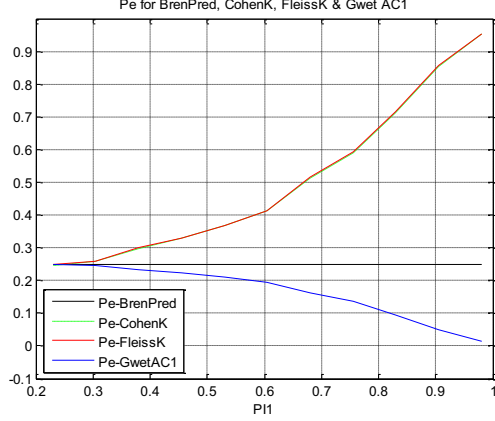


Figure 2: Graph of  $P_e$  against  $\pi_l$  for unweighted Cohen and Fleiss Kappas and the Brennan-Prediger and Gwet  $AC_1$  coefficients.

### IX. GWET'S $AC_1$ AND $AC_2$ COEFFICIENTS

Gwet's  $AC_1$  and  $AC_2$  coefficients can be considered generalisations of the unweighted and weighted Brennan-Prediger coefficients and they use the same value of unweighted or weighted  $P_o$ . In order to calculate a value of  $P_e$  (probability of agreement by chance) Gwet [3] uses the idea, similar to that used by Aickin [12], of dividing the subjects into those which are 'hard to score' and those which are 'easy to score', and estimating  $P_e$  for the 'hard' subjects only. It is suggested that the 'easy' subjects may be disregarded on the grounds that any agreement for easy subjects will not be by chance. Gwet implements this idea in a probabilistic way by defining a function  $P(R)$  as the probability of selecting a subject that is hard to score.  $P(R)$  is estimated from the distribution of the scores given to all subjects by the  $R$  raters. The degree to which all raters give the same or similar scores to the  $N$  subjects is presumed to determine the degree to which the subjects are easy to score. A group of subjects whose scores are more evenly distributed over the  $Q$  available scores is presumed to be harder to score, since there is less agreement in the scoring. Gwet's formula for  $P(R)$  is equation (18):

$$P(R) = \frac{\sum_{k=1}^Q \pi_k (1 - \pi_k)}{1 - 1/Q} \quad (18)$$

The upper curve in Figure 3 shows  $P(R)$  plotted against  $\pi_l$  for  $N=50$ ,  $R=5$  and  $Q=4$ . When  $\pi_l = 0.25$  the distribution of scores is even and all scores are equally likely. In this case, equation (18) gives a probability  $P(R)$  of selecting a hard subject close to 1. As  $\pi_l$  is increased towards 1,  $P(R)$  decreases towards zero.

The formulae for Gwet's  $AC_1$  and  $AC_2$  coefficients are obtained by modifying the Brennan-Prediger coefficient as follows:

$$AC_1 = \frac{P_o - P_e}{1 - P_e} \quad \text{where} \quad P_e = \frac{P(R)}{Q} \quad (19)$$

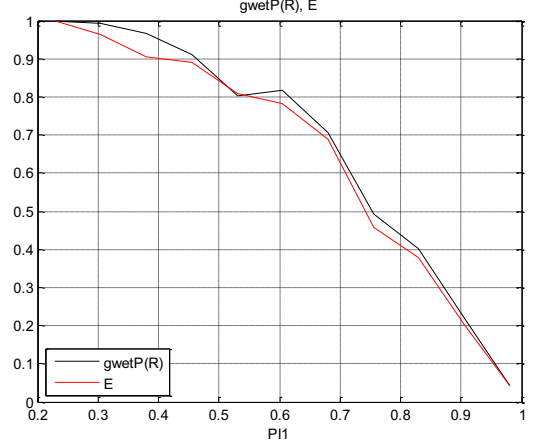


Figure 3. Gwet's  $P(R)$  function and function  $E$

$$AC_2 = \frac{P_o - P_e}{1 - P_e} \quad \text{where} \quad P_e = \frac{1}{Q^2} \sum_{k=1}^Q \sum_{\ell=1}^Q w(k, \ell) P(R) \quad (20)$$

$AC_1$  applies to categorical or unweighted scoring and  $AC_2$  to ordinal scoring with weighting. The effect of  $P(R)$  is to reduce the estimation of  $P_e$  as the probability of subjects being hard to score decreases.  $AC_1$  is plotted against  $\pi_l$  in Figure 1, and the corresponding values of  $P_e$  for  $AC_1$  are plotted in Figure 2.

As may be expected, the behaviour of  $AC_1$  is similar to that of the Brennan-Prediger coefficient. The effect of  $P(R)$  is to decrease  $P_e$  as the degree of concentration on a single score increases, and thus to increase  $AC_1$  in comparison to the Brennan-Prediger coefficient which continues to assume an even distribution of scores.

It seems reasonable to assume that  $P_e$  should decrease as the distribution of scores decreases. In this case, the subjects are indeed likely to become easier to score. But there are other cases where subjects may become easier to score, but the scores are not concentrated on a single score.

Consider, for example, the scores given in Table 2. Almost all scores agree, therefore it must be inferred that these subjects are easy to score. But  $P(R)$  as defined by Gwet [3] will be close to 1 for this example, and rightly so. The description of  $P(R)$  as the probability of selecting a subject that is hard to score is therefore misleading. It is better to describe  $P(R)$  in terms of the degree to which the sample of  $N$  subjects is likely representative of a typical population of subjects. Despite the explanation given by Gwet [3],  $P(R)$  is really defined from the overall distribution of scores and not the hardness or easiness of scoring.

### X. APPLICATION OF GWET'S $P(R)$ FORMULA TO COHEN AND FLEISS KAPPAS

Gwet [3] states that it would not be appropriate to take marginal probabilities (i.e. score distributions) into account when defining  $AC_1$  and  $AC_2$ . Despite this assertion, there seems a case for applying a measure of the degree to which the  $N$  subjects are representative of the population to both the Cohen and Fleiss Kappas. Taking  $P(R)$  as such a measure,

multiplying equations (11) and (16) for the Cohen and Fleiss Kappas respectively gives the graphs referred to as CohenK-AC1 and FleissK-AC1 in Figure 4. The graphs almost coincide, and the paradox exhibited by the Cohen and Fleiss Kappas has now been eliminated.

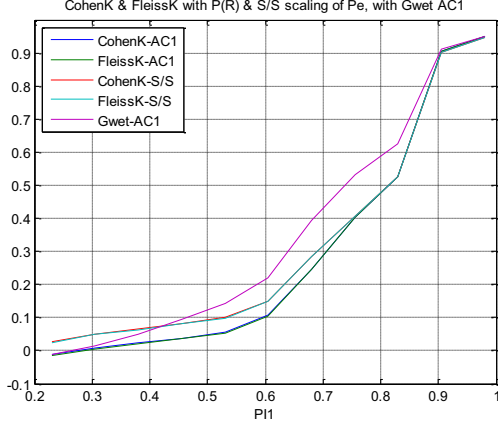


Figure 4. Cohen and Fleiss Kappas with  $P(R)$  and subject-by-subject scaling compared with Gwet's  $AC_1$  coefficient.

#### XI. SUBJECT-BY-SUBJECT IMPLEMENTATION OF THE GWET PRINCIPLE

A more direct way of implementing the principle discussed by Gwet [3] is to apply it to each individual subject rather than applying it probabilistically. For each subject  $i$ , define a 'by chance' probability,  $E(i)$ , according to the rater scores it has been given and the overall spread of rater scores. Then the contribution to  $P_e$  of any 'by chance' disagreement within rater pairs scoring subjects  $i$  and  $j$  may be scaled according to  $E(i)$  and  $E(j)$ .

For Cohen Kappa,  $P_e$  as previously defined by equation (11) becomes:

$$P_e = \frac{1}{LN^2} \sum_{i=1}^N \sum_{j=1}^N \sum_{r=1}^R \sum_{s=r+1}^R w(A(i,r), A(j,s)) \max(E(i), E(j)) \quad (21)$$

The contribution to  $P_e$  from scores  $A(i,r)$  and  $A(j,s)$  are unaffected if either subject  $i$  or subject  $j$  is considered likely to have been scored by chance. If both subjects are considered less likely to have been scored by chance, the maximum of  $E(i)$  and  $E(j)$  will be reduced, perhaps to zero. Therefore, the contribution to  $P_e$  from scores  $A(i,r)$  and  $A(j,s)$  will be reduced, thus reducing the estimated probability of agreement by chance.

For the Brennan-Prediger coefficient,  $P_e$  in equation (9) becomes:

$$P_e = \frac{T_w}{Q^2} \sum_{i=1}^N \sum_{j=1}^N \max(E(i), E(j)) \quad (22)$$

and for the Fleiss Kappa:

$$P_e = \frac{1}{R^2 N^2} \sum_{i=1}^N \sum_{j=1}^N \sum_{r=1}^R \sum_{s=1}^R w(A(i,r), A(j,s)) \max(E(i), E(j)) \quad (23)$$

This new approach requires a definition of 'by chance' probability  $E(i)$  for each subject. There are several interesting ways to do this, but a simple and obvious one is to use equation (24) as suggested by equation (18) which defines  $P(R)$ .

$$E(i) = \frac{\sum_{k=1}^Q \pi_k(i)(1 - \pi_k(i))}{1 - 1/Q} \quad \text{with} \quad \pi_k(i) = \frac{1}{R} \sum_{r=1}^R \delta(A(i,r), k) \quad (24)$$

Note that:

$$\pi_k = \frac{1}{N} \sum_{i=1}^N \pi_k(i) \quad \text{and} \quad E = \frac{1}{N} \sum_{i=1}^N E(i) \neq P(R) \quad (25)$$

Therefore we do not expect the subject-by-subject scaling by  $E(i)$  to be equivalent to the scaling by  $P(R)$  as defined by Gwet. Differences between this subject-by-subject approach and Gwet's  $P(R)$  implementation may be illustrated by comparing the curves in Figure 3. The black curve is  $P(R)$  and the red curve represents  $E$  as defined by equation (25). Both curves in Figure 3 reduce to zero as subjects become more and more concentrated on a single score. However if scaled to a maximum of 1,  $E$  remains close to  $P(R)$ .

The probability measure defined by equation (24) was applied subject-by-subject to the unweighted Cohen and Fleiss Kappas by redefining  $P_e$  as in equations (21) and (23). The new forms of Kappa are referred to as CohenK-S/S and FleissK-S/S. The result is seen in Figure 4, and may be compared with the other measures. It may be seen that the Kappas almost coincide and are close to but generally lower than the  $AC_1$  coefficient. The modification has eliminated the paradox that kept the Fleiss and Cohen Kappas close to zero in Figure 1.

#### XII. INTRA-CLASS CORRELATION (ICC)

Measurements of consistency in ordinal scoring are forms of correlation. The Pearson Correlation coefficient [13] is not normally appropriate for measuring consistency [14] as it takes into account only variations about the mean for each rater. However, the 'intra-class correlation' coefficient (ICC) [15] may be used as a consistency measure. The original form of ICC [16] for a pair of raters A and B may be written as follows:

$$ICC = \frac{\sum_{i=1}^N (A(i) - m)(B(i) - m)}{0.5 \left( \sum_{i=1}^N (A(i) - m)^2 + \sum_{i=1}^N (B(i) - m)^2 \right)} \quad (26)$$

where  $N$  subjects are scored  $\{A(i)\}_{1,N}$  by rater  $A$  and  $\{B(i)\}_{1,N}$  by rater  $B$ , and  $m$  denotes the mean of the scores given by both  $A$  and  $B$ .

It may be shown that [2] ICC is exactly equal to quadratically weighted Fleiss Kappa, and therefore incorporates correction for chance agreement. Consequently, it produces the Gwet paradox when there is a concentration on one score. As with the Cohen and Fleiss Kappas, it may be modified by the application of Gwet's  $P(R)$  function or its subject-by-subject implementation.

### XIII. CONCLUSIONS

As reported by Gwet [3], there is a fundamental flaw with the Cohen and Fleiss Kappas and this is also manifest with  $ICC$  when used for measuring rater agreement. The flaw leads to the paradox reported by Gwet which occurs when rater scores are concentrated on one score. It arises because the subjects and rater scores provided are used not only to quantify the actual agreement but also to estimate the probability of agreement by chance ( $P_e$ ).

The basic problem is that  $P_e$  is inadequately defined since no explicit assumption is made about the known or assumed statistics of the population from which the subjects are selected. When using the Cohen or Fleiss Kappa, it must be assumed that the sample of subjects provided is representative of this population, though this is usually not stated. The paradox occurs because the subjects and their scores are considered by the user not to be representative. The results just look wrong for the paradox case, but in fact they are correct if the assumption is made that the sample is representative of the population.

If it is assumed that all scores are equally likely in the population, the Brennan-Prediger coefficient correctly estimates  $P_e$ . As a means of catering for other distributions of scores, Gwet sets out to improve this coefficient by de-emphasising the contributions to  $P_e$  from subjects that are considered easy to score. These subjects are considered unlikely to have been scored by chance. The de-emphasis is achieved by multiplying each contribution by a function  $P(R)$  which is dependent on the distribution of scores.  $P(R)$  is described by Gwet as the probability of a subject within the sample being hard to score. But this description is misleading and the function would be better described in terms of the degree to which the distribution of scores is representative of an assumed population of subjects. Despite the misleading description, there is a case for applying  $P(R)$  to the Cohen and Fleiss Kappas and  $ICC$  to eliminate the paradox that may be exhibited by these coefficients.

We have investigated a subject-by-subject implementation of the Gwet principle with a 'by chance' probability estimated for each individual subject. It eliminates the paradox and has potential to improve the underlying estimate of the population statistics from the sample provided.

### REFERENCES

- [1] Gadepalli C., Jalalinajafabadi F, Xie Z, Cheetham BMG & Homer JJ, "Voice Quality Assessment by Simulating GRBAS Scoring," in proceeding of UKSim-AMSS 11<sup>th</sup> European Modelling Symposium on Mathematical Modelling and Computer Simulation (EMS2017), Manchester, UK, November 2017.
- [2] Xie Z., Gadepalli C. & Cheetham BMG, "Reformulation and Generalisation of the Cohen and Fleiss Kappas," *LIFE: International Journal of Health and Life-Sciences*, Vol 3 no 2, November 2017.
- [3] Gwet K. L., Handbook of inter-rater reliability: The definitive guide to measuring the extent of agreement among raters, Advanced Analytics, LLC, 2014.
- [4] Brennan RL & Prediger DJ, Coefficient Kappa: some uses misuses and alternatives, *Educational and Psychological Measurement* 41, 1981, pp.687-699.
- [5] Cohen J., A coefficient of agreement for nominal scales, *Educational and Psychosocial Measurement*, 20(1), 1960, pp.37-46.
- [6] Hubert I., Kappa Revisited, *Psychol Bull*, 84, 1977, pp.289-297.
- [7] Conger A.J., Integration and Generalisation of Kappas for Multiple Raters, *Psychol Bull.*, 88, 1980, pp.322-328.
- [8] Light R.J., Measures of response agreement for qualitative data: some generalisations and alternatives, *Psychol Bull*, 76, 1971, pp.365-377.
- [9] Cohen J., Weighted Kappa: Nominal scale agreement provision for scaled disagreement or partial credit, *Psychological Bulletin*, 70(4), 213,1968.
- [10] Fleiss J.L., Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76 (5), 1971, pp.378-382.
- [11] Scott WA, Reliability of Content Analysis: the case of nominal scale coding, *Public Opinion Quarterly* XIX, 1955, pp. 321-325
- [12] Aickin M., Maximum Likelihood Estimation of Agreement in the Constant Predictive Probability Model and its Relation to Cohen's Kappa, *Biometrics*, 46, 1990, pp. 293-302
- [13] Lee Rodgers J. & Nicewander W.A., Thirteen ways to look at the correlation coefficient, *The American Statistician*. 42(1), 1998, pp. 59-66.
- [14] Bland J.M. & Altman D., Statistical methods for assessing agreement between two methods of clinical measurement, *The Lancet*, 327(8476), 1986, pp.307-310.
- [15] Koch G.G., Intraclass correlation coefficient. *Encyclopedia of Statistical Sciences*, 1982.
- [16] Rödel E., Fisher R.A., *Statistical Methods for Research Workers*, 14. Aufl., Oliver & Boyd, Edinburgh, London. XIII, 362 S., 12 Abb., 74 Tab., 40 s. *Biometrical Journal*. 13(6), 1971, pp.429-30.

TABLE I. DISTRIBUTION OF SUBJECT SCORES ILLUSTRATING THE GWET PARADOX

Rater	Scores for subjects 1-20																		
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
2	1	1	2	1	1	1	1	1	1	1	1	1	1	1	1	1	2	1	1

TABLE II. DISTRIBUTION OF SCORES FOR SUBJECTS LIKELY TO BE EASY TO SCORE

Rater	Scores for subjects 1-20																		
1	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1
2	1	2	2	2	1	2	1	2	1	2	1	2	1	2	1	1	1	2	1