

1 **Explaining the unsuitability of the kappa coefficient in the assessment and comparison of the**  
2 **accuracy of thematic maps obtained by image classification**

3

4 Giles M. Foody

5 School of Geography, University of Nottingham, Nottingham, NG7 2RD, UK

6

7 Key words: accuracy, kappa coefficient, chance, prevalence, bias.

8

9 **Abstract**

10 The kappa coefficient is not an index of accuracy, indeed it is not an index of overall agreement but  
11 one of agreement beyond chance. Chance agreement is, however, irrelevant in an accuracy assessment  
12 and is anyway inappropriately modelled in the calculation of a kappa coefficient for typical remote  
13 sensing applications. The magnitude of a kappa coefficient is also difficult to interpret. Values that  
14 span the full range of widely used interpretation scales, indicating a level of agreement that equates to  
15 that estimated to arise from chance alone all the way through to almost perfect agreement, can be  
16 obtained from classifications that satisfy demanding accuracy targets (e.g. for a classification with  
17 overall accuracy of 95% the range of possible values of the kappa coefficient is -0.026 to 0.900).  
18 Comparisons of kappa coefficients are particularly challenging if the classes vary in their abundance  
19 (i.e. prevalence) as the magnitude of a kappa coefficient reflects not only agreement in labelling but  
20 also properties of the populations under study. It is shown that all of the arguments put forward for the  
21 use of the kappa coefficient in accuracy assessment are flawed and/or irrelevant as they apply equally  
22 to other, sometimes easier to calculate, measures of accuracy. Calls for the kappa coefficient to be  
23 abandoned from accuracy assessments should finally be heeded and researchers are encouraged to  
24 provide a set of simple measures and associated outputs such as estimates of per-class accuracy and  
25 the confusion matrix when assessing and comparing classification accuracy.

26 **1. Introduction**

27 The kappa coefficient of agreement was introduced to the remote sensing community in the early  
28 1980s as an index to express the accuracy of an image classification used to produce a thematic map  
29 (Congalton et al., 1983; Rosenfield and Fitzpatrick-Lins, 1986). Early papers highlighted the  
30 limitations of conventional approaches to accuracy assessment, especially the omnibus index of  
31 overall accuracy that indicates the proportion of correctly classified cases (Turk, 1979). A major  
32 concern with the latter is that its magnitude can be highly sensitive to variations in class abundance  
33 (i.e. it is prevalence dependent). This problem can be easily illustrated in relation to a basic binary  
34 classification such as that used in studies of land cover change. If one class is very rare, as change  
35 typically is, an apparently very accurate classification could be achieved by simply allocating all cases  
36 to the most abundant class (Fielding and Bell, 1997; Hoehler, 2000). In such circumstances the overall  
37 accuracy would seem to be very high but the map produced with the classification would actually  
38 provide a very poor representation of the classes, especially with regard to the rare class that may be  
39 of particular interest.

40

41 To address the problems associated with overall accuracy, the community has been encouraged to  
42 estimate and communicate with it measures of per-class accuracy (Story and Congalton, 1986;  
43 Janssen and van der Wel, 1994; Congalton and Green, 2009; Stehman and Foody, 2009; Olofsson et  
44 al., 2014) as well as explore other measures of accuracy and its reporting (e.g. Finn, 1993; Pontius,  
45 2000; Liu et al., 2007; Foody, 2011; Pontius and Millones, 2011; Comber et al., 2012; Pontius and  
46 Parmentier, 2014; Tsutsumida and Comber, 2015; Ye et al., 2018; Ariza-Lopez et al., 2019). For  
47 example, the conditional probability that a case has been allocated a class label that corresponds to its  
48 actual class of membership which is often referred to as producer's accuracy (Congalton and Green,  
49 2009; Stehman and Foody, 2009; Olofsson et al., 2014) can indicate accuracy on a per-class basis.  
50 Similarly, per-class accuracy could be assessed by relating the number of correctly classified cases of  
51 a class to the number of cases allocated to that class in the classification and this is often referred to as  
52 user's accuracy (Congalton and Green, 2009; Stehman and Foody, 2009; Olofsson et al., 2014). The

53 desire for a single omnibus measure, however, encouraged the exploration of measures of accuracy  
54 that seek to summarise accuracy over all classes in a single index and address impacts of issues such  
55 as class abundance on the apparent accuracy. Indeed the kappa coefficient was proposed as an index  
56 that improved upon overall accuracy (Ubersax, 1987; Maclure and Willett, 1987) and in the remote  
57 sensing community it has been promoted as being an advancement on overall accuracy (Congalton et  
58 al., 1983; Fitzgerald and Lees, 1994).

59

60 Key arguments put forward for the adoption of the kappa coefficient as an index of classification  
61 accuracy were along the lines that it corrected for chance agreement, scales exist for its interpretation,  
62 it may be estimated on a per-class as well as on an overall basis and that a variance term may be  
63 estimated for it allowing statistically rigorous comparisons to be undertaken (Congalton et al., 1983;  
64 Rosenfield and Fitzpatrick-Lins, 1986). Perhaps because of the correction for chance agreement, it is  
65 also sometimes claimed that the kappa coefficient is relatively independent of variations in class  
66 prevalence (Manel et al., 2001).

67

68 The papers that introduced the kappa coefficient for accuracy assessment in remote sensing have had  
69 an enormous impact on the research community. These papers have been very highly cited and have  
70 been followed by other hugely influential publications that have further promoted the use of the kappa  
71 coefficient in accuracy assessment (e.g. Congalton, 1991; Congalton and Green, 2009). These  
72 publications have helped to foster the widespread use of the kappa coefficient that has been aided by  
73 the inclusion of functionality for its calculation in popular image processing software (Pontius and  
74 Millones, 2011).

75

76 Despite the widespread promotion of the kappa coefficient and the ease of its estimation, there are  
77 many concerns with its use in accuracy assessment. Although widely used, the kappa coefficient has  
78 had a troubled history, with concerns ranging from the use of incorrect equations (Fleiss et al., 1969;

79 Rosenfield and Fitzpatrick-Lins, 1986; Hudson and Ramm, 1987) to more fundamental calls for the  
80 kappa coefficient to be abandoned (e.g. Pontius and Millones, 2011). Indeed the use of the kappa  
81 coefficient is regarded explicitly as poor practice in accuracy assessment (Olofsson et al. 2013; 2014).  
82 Sadly the calls to abandon the use of the kappa coefficient in accuracy assessment seem to have fallen  
83 on deaf ears. It may be that the kappa coefficient is still widely used because it has become ingrained  
84 in practice and there may be a sense of obligation to use it (Stehman and Foody, 2019). Indeed many  
85 researchers seem to use it because precedent for its use exists but given the concerns with the kappa  
86 coefficient this is merely an argument to allow mistakes to be repeated. Mistakes happen, but should  
87 be used as a positive learning experience that leads to constructive change rather than a situation to be  
88 repeated.

89

90 It is unclear why the calls to abandon the use of the kappa coefficient in accuracy assessment have not  
91 been heeded as the criticisms have been damning with recommendations for good practice clear (e.g.  
92 Foody, 1992; Stehman, 1997a; Pontius and Millones, 2011; Stehman and Foody, 2009; Olofsson et  
93 al., 2013, 2014). It may be that theoretical arguments have been challenging or that the ease with  
94 which the kappa coefficient may be estimated as relevant functionality is often embedded in popular  
95 software leads to widespread and possibly unquestioning use. For example, in the period after the  
96 publication of the ‘death to kappa’ paper by Pontius and Millones (2011), the kappa coefficient was  
97 reported in half of the relevant literature (Morales-Barquero et al., 2019). The use of the kappa  
98 coefficient seems to be embedded into standard practice despite well-known concerns that have been  
99 widely disseminated. One possible reason for this unsatisfactory situation is that the community is  
100 unaware of the magnitude of the problems associated with the use of the kappa coefficient. Hence,  
101 this article aims to revisit major concerns with the use of the kappa coefficient to demonstrate its  
102 unsuitability as an index of classification accuracy in remote sensing using simple examples with a  
103 focus on highlighting the challenges of interpreting a kappa coefficient by stressing the difficulties in  
104 interpreting its magnitude. It will be stressed that all of the arguments put forward for the use of the  
105 kappa coefficient are flawed or, in the sense that they are not unusual or unique, irrelevant. The article

106 will first review the estimation of the kappa coefficient and key attributes that have been espoused in  
107 support of its use. The latter will be critically evaluated to highlight key concerns before providing  
108 some simple examples to demonstrate the problems that can be encountered in the interpretation of  
109 the magnitude of a kappa coefficient. Throughout the focus is on commonly encountered situations  
110 and hence limited to evaluations of standard hard classifications.

111

## 112 **2. Estimation of the kappa coefficient**

113 The kappa coefficient can be estimated easily from the confusion or error matrix that is widely used in  
114 classification accuracy assessment. For ease of discussion, the main focus will be on the simplest case  
115 of a binary confusion matrix which is widely used in, for example, studies of land cover change  
116 (Figure 1). The approach readily extends to larger, multi-class, matrices and this is briefly discussed  
117 for completeness. For ease of presentation, it will also be assumed throughout that the sample of cases  
118 used to form the confusion matrix was acquired using simple random sampling unless stated  
119 otherwise; different sampling designs can be used and the correct formulae for use with them are  
120 provided in the literature (e.g. Stehman, 1996, 1997b).

121

122 In a binary classification there are just two classes. Thus, in the map produced by a binary image  
123 classification, each case (e.g. image pixel) either has (+) a particular trait associated with it or it has  
124 not (-). For example, in a remote sensing application the case might be labelled in the map as  
125 representing an area of change or of no change. Similarly, the labels might be forest and non-forest or  
126 urban and non-urban or to some other specific class of interest or not. Critically, a case may also have  
127 similar labels applied to it in a ground reference data set used to assess classification accuracy. The  
128 cross-tabulation of the class labels observed in the map and those in the reference data set yields a  
129 basic 2 x 2 confusion matrix, often referred to as an error matrix, from which a range of summary  
130 measures of classification accuracy can be obtained (Figure 1). Based on the assumption that the map

131 and reference data sources are considered to be two independent raters, the kappa coefficient of  
 132 agreement may be estimated from this matrix.

		Reference ↓			
		Class	+	-	Σ
Map ↓	+	<i>a</i>	<i>b</i>	<i>n</i> <sub>+</sub>	
	-	<i>c</i>	<i>d</i>	<i>n</i> <sub>-</sub>	
	Σ	<i>n</i> <sub>+</sub>	<i>n</i> <sub>-</sub>	<i>n</i>	

133  
 134

135 Figure 1. The confusion matrix for a binary classification based on a simple random sample of *n*  
 136 cases.

137

138 Before exploring the estimation of the kappa coefficient further it may be useful to focus first on the  
 139 composition of the confusion matrix. The binary confusion matrix has four elements that summarise  
 140 every possible scenario of class labelling. The number of cases with each of the four possible class  
 141 allocation scenarios, *a-d*, are inserted into the appropriate matrix elements. Of these, *a* cases are  
 142 labelled as having the trait of interest in both the image classification that forms a thematic map and  
 143 the reference data; these are often termed true positives. The *d* cases that are labelled as not having the  
 144 trait of interest in both the image classification and the reference data lie in the other element of the  
 145 matrix's main diagonal; these are often termed true negatives. Thus, the cases lying in elements of the  
 146 main diagonal, *a* and *d*, represent those that have been correctly classified. All of the cases that have  
 147 been incorrectly classified lie in the off-diagonal elements of the matrix. Of these, *b* are those cases  
 148 that have been classed as having the trait of interest but do not actually possess it; these are commonly  
 149 referred to as false positives. Such cases represent commission errors, sometimes referred to as type I  
 150 errors although the use of this terminology can sometimes be problematic (Thron and Miller, 2015).  
 151 Finally, *c* cases have the trait of interest in the reference data but were classified as not having it; these  
 152 are commonly referred to as false negatives. These latter cases represent omission errors, sometimes  
 153 referred to as type II errors. The cases on which the classification and reference data differ in labelling  
 154 are the misclassifications or errors. In Figure 1, omission is assessed with a focus on the columns of  
 155 the matrix while commission is assessed with a focus on the rows of the matrix. The total number of

156 cases lying in each row and each column can be determined by summing the relevant matrix elements.  
157 These row and column total values are often referred to as the matrix marginal values. Their total,  
158 calculated over all rows or all columns, also equates to the total number of cases,  $n$ , used to form the  
159 matrix. The difference between the row and column proportions for a class indicate non-site specific  
160 accuracy and indicate map bias which is sometimes referred to as quantity disagreement (Pontius and  
161 Millones, 2011; Stehman and Foody, 2019). Finally, the prevalence,  $\theta$ , of the trait of interest which  
162 indicates its abundance may be estimated from  $\frac{(a+c)}{n} = \frac{n_{.+}}{n}$  and is a property of population being  
163 studied. Ideally, a measure of accuracy should reflect only the quality of the classification and not  
164 vary with prevalence. Indeed, the prevalence dependency of overall accuracy noted at the beginning  
165 of this article is one of its major limitations as a measure of accuracy. Some measures, such as  
166 producer's accuracy, are prevalent independent if the diagnostic ability of the classifier is unaffected  
167 by prevalence, which can aid their interpretation; in common remote sensing applications the  
168 producer's accuracy may, however, be expected to be prevalent dependent.

169

170 Using notation similar to Cohen (1960), the kappa coefficient of agreement,  $\kappa$ , is estimated from:

171 
$$\kappa = \frac{p_o - p_e}{1 - p_e} \quad (1)$$

172 where  $p_o$  is the proportion of cases correctly classified (i.e. overall accuracy) and  $p_e$  is the expected  
173 proportion of cases correctly classified by chance; note with this notation the distinction between  
174 parameters and estimated parameters is not explicit but the text will indicate where sample-based  
175 estimates are being made or used. The magnitude of  $\kappa$  lies on a scale from -1 to +1 but interest is  
176 typically focused on only on positive values because negative values indicate a level of agreement less  
177 than that due to chance and can be difficult to interpret (Sim and Wright, 2005). The maximum value  
178 of +1 occurs when there is perfect agreement and a value of 0 arises when the observed agreement  
179 equals that due to chance (Cohen, 1960). Commonly the magnitude of the kappa coefficient is  
180 interpreted relative to a scale. One such interpretation scale that has been widely used in remote  
181 sensing applications is that proposed by Landis and Koch (1977).

182 Central to the estimation of the kappa coefficient is the estimation of the level of agreement and also  
183 the level of agreement that occurs due to chance. For the simple case of a binary confusion matrix  
184 such as shown in Figure 1, the proportion of agreement,  $p_o$  is estimated from

$$185 \quad p_o = \frac{a + d}{n} \quad (2)$$

186 in which  $a$  and  $d$  are the number of cases correctly labelled (i.e. the true positive and true negative  
187 cases), lying in the elements of the main diagonal of the confusion matrix (Cohen, 1960; Congalton et  
188 al., 1983). Thus,  $p_o$  is simply the sum of all correctly classified cases divided by the total number of  
189 cases used to form the matrix and expresses the proportion of correctly labelled cases (i.e. overall  
190 accuracy); it is often multiplied by 100 and expressed as a percentage which is commonly termed the  
191 percentage correctly classified cases. Although an imperfect index of accuracy, the proportion of  
192 correctly allocated cases is relatively easy to estimate and understand (Pontius and Millones, 2011).  
193 Before going into any further detail one thing to note at this stage of the discussion is that the kappa  
194 coefficient is estimated from  $p_o$ , it is an additional analytical step required after the estimation of  
195 overall accuracy.

196

197 There are a variety of ways to estimate chance agreement (Byrt et al., 1993), but the version that is  
198 adopted commonly in remote sensing, which is used in the estimation of Cohen's kappa coefficient, is  
199 based on a simple analysis of the row and column marginal values (Byrt et al., 1993; Lantz and  
200 Nebenzahl, 1996; Hoehler, 2000, Sim and Wright, 2005). In this, the proportion of agreement  
201 expected due to chance,  $p_e$ , may be obtained from equation 3.

$$202 \quad p_e = \left( \left( \frac{a+c}{n} \right) \left( \frac{a+b}{n} \right) \right) + \left( \left( \frac{b+d}{n} \right) \left( \frac{c+d}{n} \right) \right) \quad (3)$$

203 Chance may be modelled differently yielding alternatives to equation 3 and these may be used in  
204 equation 1 to yield other indices of agreements. For example, Scott's pi,  $\pi$ , is estimated from equation  
205 1 but, as it is based on different assumptions to the kappa coefficient, the estimation of  $p_e$  is different  
206 (Byrt et al., 1993; Banerjee et al., 1999).

207

208 To illustrate accuracy on a per-class basis it is possible to estimate the conditional kappa coefficient  
209 (Rosenfield and Fitzpatrick-Lins, 1986; Czaplewski, 1994; Congalton and Green, 2009). For the class  
210  $i$ , which has either the + or – label, the latter may be estimated from

$$211 \quad \kappa_i = \frac{nn_{ii} - n_i \cdot n_{\cdot i}}{nn_{i\cdot} - n_i \cdot n_{\cdot i}} \quad (4)$$

212 The variance for kappa may be estimated (Congalton et al., 1983; Congalton and Green, 2009) and  
213 can be usefully expressed in terms of the standard error,  $\sigma_\kappa$ , which is the square root of the variance.  
214 The details of the estimation are not central to the argument in this article but the equation for its  
215 estimation for those interested is given in Figure 2. A large literature discusses the estimation of the  
216 variance and related terms in more detail (e.g. Fleiss et al., 1969, 2013; Hudson and Ramm, 1987;  
217 Czaplewski, 1994).

218

219 The standard error may be used to define confidence limits around the estimated value of a kappa  
220 coefficient. For example, the 95% confidence interval (95% CI) would be  $\kappa \pm 1.96\sigma_\kappa$  as at this level  
221 of confidence the standard score,  $z$ , is 1.96. The statistical significance of a kappa coefficient may also  
222 be assessed, using:

$$223 \quad z = \frac{\kappa}{\sigma_\kappa} \quad (5)$$

224 which indicates the degree to which the level of agreement observed is better than that arising from  
225 chance alone (Congalton and Green, 2009; Fleiss et al., 2013). More usefully, this also provides the  
226 basis to compare an estimated kappa coefficient against other values and also to compare the  
227 difference between two estimated kappa coefficients. This is particularly useful when seeking to  
228 undertake a statistically rigorous and credible comparison of the accuracy of two thematic maps. For  
229 example, two maps, A and B, may have been produced for a region using two different classifiers and  
230 the researcher may be interested in knowing if they differ in accuracy. The test for the significance of  
231 the difference between two kappa coefficients estimated using independent samples is:

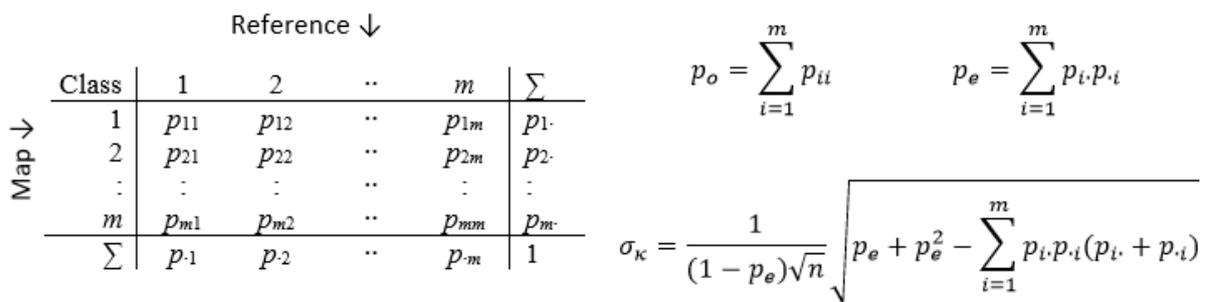
232 
$$Z = \frac{\kappa_A - \kappa_B}{\sqrt{\sigma_{\kappa A}^2 + \sigma_{\kappa B}^2}} \quad (6)$$

233 where  $\kappa_A$  and  $\kappa_B$  are the estimated kappa coefficients for maps A and B respectively, and  $\sigma_{\kappa A}$  and  $\sigma_{\kappa B}$   
 234 are the associated estimates of the standard error of kappa for maps A and B respectively (Cohen,  
 235 1960; Congalton and Mead, 1983; Congalton et al., 1983; Rosenfield and Fitzpatrick-Lins, 1986,  
 236 Smits et al., 1999). Two maps would be deemed to be of different accuracy if  $|z| > 1.96$  at the 95%  
 237 level of confidence. If the hypothesis under test has a directional component (e.g. that one map is  
 238 more accurate than another) a one-sided rather than two-sided test can be undertaken in the usual way  
 239 (Foody, 2009; Fleiss et al., 2013).

240

241 The discussion in this article is focused on binary classifications for ease but the issues extend to  
 242 multi-class classifications. For multi-class classifications the nature of the confusion matrix and key  
 243 equations are given in Figure 2.

244



245

246

247 Figure 2. The confusion matrix for a multi-class classification involving  $m$  classes, expressed as  
 248 proportions, together with key equations for the estimation of the kappa coefficient and its standard  
 249 error.

250

251

### 252 3. Challenging the arguments for the use of the kappa coefficient

253 Before addressing the substantive problems with the kappa coefficient it should be noted that a range  
254 of problems have been encountered in its use in remote sensing. For example, there is often a failure  
255 to recognise impacts of the sample design used to acquire the cases used in estimation (Stehman,  
256 1996), incorrect variance equations have been used (Rosenfield and Fitzpatrick-Lins, 1986), and many  
257 comparative assessments have used related rather than independent samples (Foody, 2004) or not  
258 recognised the directionality of the study which may require testing for dissimilarities related to  
259 inferiority, superiority or equivalence rather than just a difference (Foody, 2009). Similar concerns  
260 could be flagged in relation to other indices of accuracy and so such problems are not the central issue  
261 of concern to this article. Here, the concern is that the kappa coefficient is unsuitable for use in  
262 accuracy assessment, the additional problems encountered in practical application are of very  
263 secondary importance. Consequently, the latter are not discussed further especially as such  
264 methodological errors are often easy to address with, for example, equations for use with stratified  
265 samples (Stehman, 1996) and cluster samples (Stehman, 1997b) as well as statistical tests for related  
266 samples (Donner et al., 2000; Foody 2004; 2009; Fleiss et al., 2013).

267

268 Central to this article are fundamental problems with the use of the kappa coefficient as an index of  
269 classification accuracy. A variety of arguments can be raised against the use of the kappa coefficient  
270 in accuracy assessment. These range from the fundamental issue that as a measure of inter-rater  
271 agreement it is not a measure of accuracy (Nishii and Tanaka, 1999; Vach, 2005; Wu et al., 2007) to  
272 substantial difficulties in its interpretation (Byrt et al., 1993; Lantz and Nebenzahl, 1996; Sim and  
273 Wright, 2005; Pontius and Millones, 2011). Here, the central focus is directed at challenging each of  
274 the arguments that have been put forward to promote the use of the kappa coefficient in order to  
275 highlight its unsuitability as a measure of classification accuracy, summarised in Table 1.

276

277 Table 1. A summary of the seven main arguments offered for the adoption of the kappa coefficient  
 278 and a brief critique of each, highlighting the argument to be either seriously flawed or irrelevant, in  
 279 the sense that while it may be a valid statement there is nothing unusual or different to other standard,  
 280 often simpler, indices of accuracy. In short, not a single one of the key arguments put forward for the  
 281 use of kappa has any real merit, each is either deeply flawed or equally applicable to other indices.

<b>Arguments for the use of kappa</b>	<b>Reality</b>
It 'corrects' for chance agreement	Flawed argument. There is no need to 'correct' for chance agreement. The source of error is unimportant in the assessment of classification or map accuracy. Furthermore, chance is an artificial construct and the way it is modelled in the estimation of $\kappa$ is inappropriate.
Its estimation is based on the entire confusion matrix	Flawed argument, indeed one that is completely untrue. The estimation is actually based on the main diagonal together with the row and column marginal totals.
It can be estimated on an overall and per-class basis	Irrelevant as the exact same can be argued for other standard measures of accuracy such as overall accuracy (i.e. the proportion of cases correctly classified) with per-class statements from the user's and producer's perspectives.
It is, to a large degree, prevalent independent	Flawed argument as untrue. Kappa is, like many other indices, very dependent on class prevalence.
A variance term may be estimated for it.	Irrelevant as the exact same can be argued for other standard measures of accuracy such as the proportion of cases correctly classified.
It allows rigorous comparison of estimates of classification accuracy.	Irrelevant as the exact same approach to comparison, which requires variance estimates, can be used with other measures of accuracy. The commonly promoted approach is also suitable for situations in which independent samples are used but often the same sample is used; methods for the comparison of accuracy estimates obtained from the same sample are available. The comparison of kappa coefficients is also problematic if there are differences in prevalence.
Scales exist for its interpretation	Flawed argument. A variety of scales exist but any scale is arbitrary and cannot be expected to be of universal applicability. The scales also ignore problems linked to issues such as class prevalence.

282

283 The kappa coefficient is designed for application to data arising from two independent raters and  
 284 provides a measure of the degree to which they agree in labelling. Indeed, an early article introducing

285 the kappa coefficient to the remote sensing community focused on its use as a measure of inter-rater  
286 agreement (Congalton and Mead, 1983). However, this type of analysis is not the scenario  
287 encountered in the assessment of classification accuracy, notably because the ground reference data  
288 are supposed to represent the true condition and the desire is to yield a measure of accuracy not  
289 simply agreement.

290

291 Classification accuracy is a measure of the quality with which a set of cases have been labelled.  
292 Fundamentally, the concern in accuracy assessment is with the amount of error or mis-labelling that  
293 has occurred in the classification. In this way the accuracy assessment is useful in terms of assessing  
294 the fitness for purpose of the classification. The latter would typically require a comparison of the  
295 estimated accuracy relative to some target value that indicates the minimum acceptable accuracy for  
296 the proposed use of the classification. A target accuracy should ideally be defined before the  
297 classification is undertaken and be tailored to the specific purpose of the classification (Foody, 2008).  
298 For example, in the pioneering work linked to Anderson (1971) and Anderson et al. (1976) for the  
299 mapping of broad land cover classes over a large area, a target of 85% correct allocation with the  
300 classes mapped to approximately equal accuracy was used. This target value was well-justified for the  
301 specific application and data sets used. For a different mapping application, a target for the specific  
302 needs of that individual application should be defined and used; the 85% target put forward by  
303 Anderson et al. (1976) is not a universally applicable one. For example, a simple binary classification  
304 involves fewer classes than the application Anderson et al. (1976) addressed and a higher target  
305 accuracy might be appropriate. An example used below, for instance, sets a target that comprises an  
306 overall accuracy of 95% with the producer's accuracy for the two classes to be at least 95%. Key  
307 attractions of this sort of Anderson-type target are that a target value can be defined in advance of the  
308 classification and it may, to some extent, help to address concerns with prevalence dependency. The  
309 latter arises because the target includes the producer's accuracy for each class and this measure of  
310 accuracy is independent of prevalence if the diagnostic ability of the classifier is fixed (Rogan and  
311 Gladen, 1978; Maclure and Willetts, 1987); but note that the valuable attribute of prevalence

312 independence is lost if the ground data set is imperfect (Foody, 2010) or if the diagnostic ability of the  
313 classifier changes with prevalence.

314

315 The desire for a target highlights an initial problem with the use of the kappa coefficient: how can a  
316 sensible target value be defined in advance of a mapping study when the marginal values of the  
317 confusion matrix are unknown? In brief, it will typically be infeasible to define a meaningful kappa  
318 coefficient as a target value in advance of the classification. It could be argued, however, that a target  
319 value is not required with the use of the kappa coefficient as the quality of the classification can be  
320 assessed relative to an interpretation scale. This will be one of the problems with the kappa coefficient  
321 that will be discussed below.

322

323 As highlighted in the introduction, several key attributes have been routinely suggested as arguments  
324 for the use of the kappa coefficient in the assessment of classification accuracy. Perhaps the most  
325 widely used argument for the adoption of the kappa coefficient is, essentially, that it corrects for  
326 chance agreement. Although the exact meaning of 'chance correction' is not always clear the core  
327 thrust appears to be that it adjusts the assessment for the effect of chance agreement; the kappa  
328 coefficient essentially quantifies the level of agreement beyond that due to chance. This is an  
329 important observation as the kappa coefficient is often treated as a measure of overall agreement  
330 rather than a measure of agreement beyond chance (Jiang and Liu, 2011) and, as noted above, chance  
331 may be modelled in different ways and so needs to be quantified with care. Because of the assessment  
332 being made relative to a random classification, which is unrealistic of real land cover mosaics, the  
333 kappa coefficient fails to meet the map relevant criterion for good practice (Stehman and Foody,  
334 2019). Moreover, the aim of an accuracy assessment is, essentially, the estimation of how much error  
335 has occurred; the lower the error the greater the accuracy. Note the origin of the error or the reason for  
336 correct labelling is of absolutely no concern to the measurement of accuracy. In a conventional  
337 accuracy assessment, a map label is either correct or it is not. There may well be interest in

338 understanding error, especially as a means to further enhance a classification-based analysis, but such  
339 assessments of skill require a different type of analysis (Turk, 1979); a distinction between the  
340 assessment of classifier performance that indicates diagnostic ability and the assessment of  
341 classification accuracy is required (Turk, 2002). Accuracy assessment merely seeks to quantify the  
342 amount of error, the origin or source of the error is irrelevant. There is, therefore, no interest in chance  
343 agreement and no desire to correct for it in a standard accuracy assessment. Indeed rather than  
344 estimate and remove the chance agreements the community should regard such agreements as a  
345 windfall gain (Turk, 2002). Even if there was a desire to explore the issue of chance agreement the  
346 estimation of its magnitude for the calculation of the kappa coefficient, equation 3, is inappropriate.  
347 Since the ground reference data represent reality rather than labels from another independent rater, it  
348 may be more appropriate to have fixed column marginal values determined by the number of classes  
349 with  $p_e = 1/m$  (Brennan and Prediger, 1981; Foody, 1992).

350

351 Another popular argument for the use of the kappa coefficient is that its variance may be estimated  
352 which facilitates rigorous testing. In particular, the ability to obtain the variance for kappa allows tests  
353 of the statistical significance of the difference between two kappa coefficients to be undertaken  
354 (Rosenfield and Fitzpatrick-Lins, 1986; Congalton and Green, 2009). These arguments are well-  
355 founded and the ability to rigorously compare estimates is a useful attribute. This situation is,  
356 however, nothing particularly special to the kappa coefficient. The variance of other estimates of  
357 accuracy such as the overall accuracy, which is simply a proportion ( $p$ ), can also be calculated. The  
358 standard error for a proportion, assuming the use of a simple random sample, can be estimated from:

359 
$$\sigma_p = \sqrt{\frac{p(1-p)}{n}} \quad (7)$$

360 Thus, the variance and related statistics can be obtained for proportions (Fleiss et al., 2013) such as  
361 overall, producer's and user's accuracy. Furthermore, contrary to claims to the reverse (Jansen and  
362 van der Wel, 1994), it is possible to rigorously compare estimates of the proportion of correctly  
363 classified cases. Thus, the statistical significance of the difference in the accuracy of two

364 classifications could be assessed using overall accuracy. The assessment would be similar to that  
365 indicated by equation (6) but with the proportion correct,  $p_o$ , and its associated variance term, which  
366 can be expressed as the standard error,  $\sigma_p$ , for each classification used instead of the kappa  
367 coefficients and their standard errors (Stehman, 1997a; Foody, 2004):

$$368 \quad z = \frac{P_{oA} - P_{oB}}{\sqrt{\sigma_{pA}^2 + \sigma_{pB}^2}} \quad (8)$$

369 Equation 8 allows the statistical significance of differences in proportions, such as overall accuracy,  
370 on the assumption that the samples used are independent. Often in remote sensing applications the  
371 same ground reference data set is used and the effect this has on the analysis could be addressed by  
372 integrating a covariance term into the test or by adopting a test suited for use with related samples  
373 such as the McNemar test as an alternative (Foody, 2004; 2009).

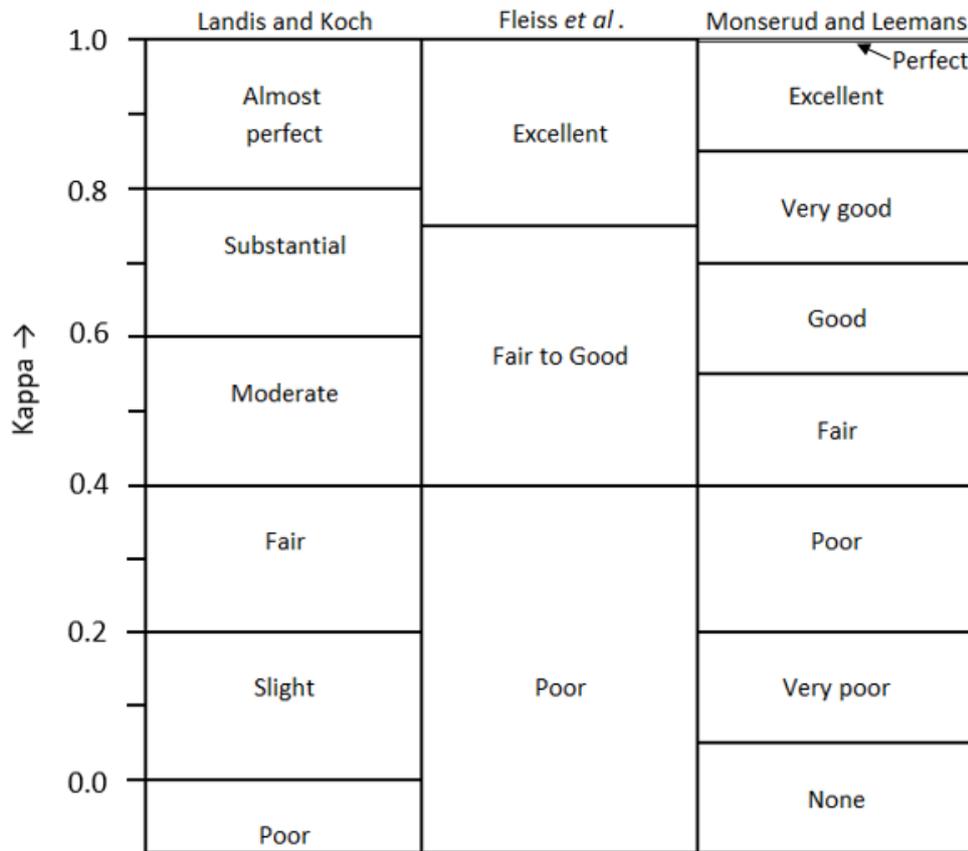
374

375 The ability to estimate a measure of accuracy on a per-class basis has also been highlighted as an  
376 advantageous feature associated with the kappa coefficient. Often referred to as conditional kappa this  
377 allows assessment on a class-specific rather than overall basis. Although this is a useful feature it is  
378 also nothing special or unique to the kappa coefficient. As noted above, per-class measures of  
379 accuracy can be obtained directly from the confusion matrix used to estimate  $p_o$ . For example, simple  
380 per-class measures such as user's and producer's accuracy can be obtained by analysing the relevant  
381 row and column of the confusion matrix depending on whether errors of commission or omission are  
382 important. For example, the producer's accuracy ( $P$ ) for the class with the trait of interest is estimated  
383 from  $P_+ = a/n_+$ ; often referred to as the true positive rate, recall or sensitivity. Similarly, the  
384 producer's accuracy may be calculated for the class without the trait of interest from  $P_- = d/n_-$ ; often  
385 referred to as specificity. Alternatively, with a focus on commission error, the user's accuracy ( $U$ )  
386 may be calculated for each class. For example, the user's accuracy for the class with the trait of  
387 interest may be estimated from  $U_+ = a/n_+$ ; often referred to as the positive predicted value or precision  
388 although this latter term should perhaps be avoided due to the potential for mis-interpretation.  
389 Sometimes researchers combine measures to yield a single summary indicator of classification

390 accuracy. One such measure which utilizes the producer's accuracy for each class is Youden's  $J$   
391 which is estimated as  $J = P_+ + P_- - 1$  (Allouche et al., 2006; Hand, 2012); sometimes referred to as the  
392 true skills statistic or informedness. This latter index is sometimes attractive as an overall summary  
393 measure of classification accuracy as the components may be prevalent independent if the diagnostic  
394 ability of the classifier is fixed and, although not without concerns, its variance may also be estimated  
395 (Allouche et al. 2006). However, there are many measures of accuracy and these can be combined in  
396 various ways. For example, average accuracy or the F1 score can be estimated. Such measures,  
397 however, are challenging to interpret and of questionable value (Stehman and Foody, 2009; Liu et al.,  
398 2007). Indeed many measures of accuracy are available and may be sensitive to different things  
399 (Hand, 2012). For a statement of map accuracy to be useful the error measure adopted should be  
400 justified and appropriate to the task in-hand (Fielding and Bell, 1997).

401

402 A key feature often used in the promotion of the use of the kappa coefficient in accuracy assessment  
403 is that scales to interpret the kappa coefficient are available. The existence of a meaningful scale could  
404 also be argued to remove the common desire for a target value in accuracy assessment. While it is true  
405 that scales for the interpretation of the kappa coefficient exist, with that provided by Landis and Koch  
406 (1977) widely used in remote sensing, there are substantial problems in their use. For example, there  
407 are a range of scales available (e.g. Figure 3) with no obvious way to choose between them and a  
408 scale could readily be constructed for other indices such as overall accuracy. More critically, it should  
409 be readily apparent that such interpretation scales are arbitrary and cannot be of universal applicability  
410 (Sim and Wright, 2005; Vach, 2005; Banerjee et al., 1999). Indeed, Landis and Koch (1977) explicitly  
411 note the arbitrary nature of the scale that they proposed in their study. Some studies may, for example,  
412 require very high quality labelling and hence the thresholds dividing the scale should be set at higher  
413 values. The arbitrary and subjective nature of the scales limit their value as a means to interpret a  
414 kappa coefficient. The problems also mean that the existence of an interpretation scale does not  
415 address the inability to define a meaningful target value if using the kappa coefficient as the index of  
416 accuracy.



417

418 Figure 3. Three scales for the interpretation of the kappa coefficient (adapted and updated from  
 419 Czaplewski, 1994). The scales are those provided by Landis and Koch (1977, page 165); Fleiss *et al.*  
 420 (2013, page 604) and Monserud and Leemans (1992, page 285). Note that the full scale of  
 421 measurement does extend to -1 but the focus is usually on positive values only.

422

423

424 The interpretation of a kappa coefficient can be challenging, especially if not accompanied by the  
 425 confusion matrix and details of the sample of cases used in its estimation. Indeed it is widely  
 426 suggested that that the provision of a kappa coefficient alone is misleading and that per-class  
 427 measures and/or indices of bias and prevalence should accompany it (Byrt *et al.*, 1993; Lantz and  
 428 Nebenzahl, 1996; Cicchetti and Feinstein, 1990); the provision of the confusion matrix and details of  
 429 the sample used in its construction would also help as they can provide the additional information  
 430 needed to interpret a kappa coefficient. A variety of challenges is encountered in interpreting the  
 431 magnitude of a kappa coefficient. In particular, two paradoxes commonly arise (Feinstein and

432 Cicchetti, 1990; Lantz and Nebenzahl, 1996; Hoehler, 2000; Sim and Wright, 2005). First, there is the  
433 situation in which there may be high level of agreement indicated by  $p_o$  but a low kappa coefficient.  
434 Second, unbalanced matrix marginal values can help produce a high kappa coefficient, especially if  
435 the marginals are asymmetrically imbalanced (Feinstein and Cichetti, 1990). These paradoxes arise  
436 because the estimation of the kappa coefficient is influenced by prevalence and bias between the  
437 raters (Byrt et al., 1993; Lantz and Nebenzahl, 1996; Hoehler, 2000). Both paradoxes can be  
438 explained by the distribution of cases within the confusion matrix. The first paradox arises because of  
439 the effect of prevalence on the estimation of the kappa coefficient and is positively related to the  
440 difference between  $a$  and  $d$  (Figure 1). The second paradox is related to bias effects that occur when  
441 the two sources of class labels used to form the confusion matrix differ in the proportion of cases with  
442 the trait of interest and varies as a function of the difference between  $b$  and  $c$  (Figure 1). Critically, the  
443 manner in which cases are distributed in the confusion matrix and its resulting marginal values can  
444 greatly impact on the magnitude of the kappa coefficient.

445

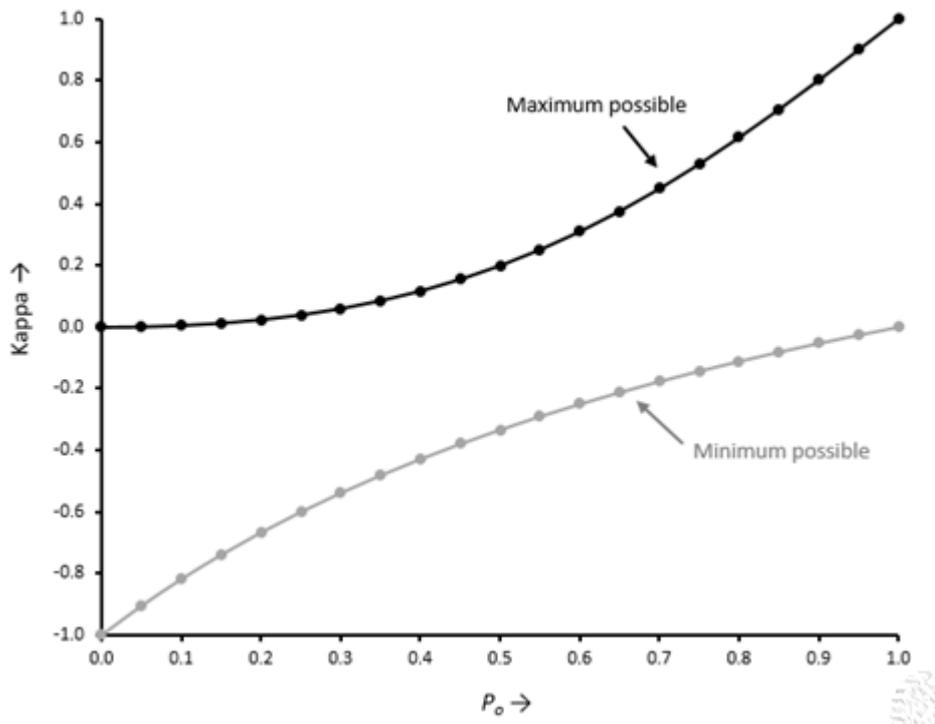
446 It is sometimes claimed that the whole confusion matrix is used in the estimation of the kappa  
447 coefficient. This claim, however, is untrue; the estimation of the kappa coefficient is based on the  
448 main diagonal and marginal values only (Nishii and Tanaka, 1999; Jiang and Liu, 2011). It is, for  
449 example, possible in a multi-class classification to change the entries in the matrix but maintain the  
450 same diagonal and marginal values and hence kappa coefficient. Because of the prevalence and bias  
451 effects noted above, knowledge of all of the elements of the matrix is, however, useful in interpreting  
452 a kappa coefficient (Lantz and Nebenzahl, 1996).

453

454 The factors that influence the magnitude of the kappa coefficient are well-known but the size and  
455 importance of the issues may not always be apparent. To help demonstrate problems in the  
456 interpretation and use of the kappa coefficient it may be helpful to explore some simple scenarios as  
457 examples. As a starting point, a range of possible values for the kappa coefficient can be obtained for

458 any given level of agreement ( $p_o$ ). This range can be explored by moving cases around the confusion  
 459 matrix in a manner that maintains the proportion of correct agreement. The maximum and minimum  
 460 kappa coefficient possible may also be estimated given an understanding of how the distribution of  
 461 cases in a confusion matrix impacts on the estimation of the kappa coefficient (Lantz and Nebanzahl,  
 462 1996). Figure 4 shows the relationship between the maximum and minimum kappa coefficient values  
 463 that can be obtained for all possible proportions of correct agreement. A key feature to note is the  
 464 extremely large difference between the maximum and minimum kappa coefficient at each value for  
 465 the proportion of correct agreement. For example, with the very high level of agreement of  $p_o=0.95$  it  
 466 would be perfectly possible for a kappa coefficient of between  $-0.026$  and  $0.900$  to be estimated.  
 467 Moreover, this very wide range of possible values for the kappa coefficient covers every single level  
 468 of the widely used interpretation scale of Landis and Koch (1977). Thus, with 95% of the cases  
 469 correctly labelled the use of the kappa coefficient could result in the level of agreement interpreted as  
 470 being anything from poor to almost perfect inclusive (Figure 3).

471



472

473 Figure 4. Relationships between the maximum and minimum possible kappa coefficient with overall

474 accuracy ( $p_o$ ).

475 The confusion matrices for the extreme values of the kappa coefficient when  $p_o=0.95$  are shown in  
 476 Figure 5 and highlight the effect of bias on the maximum value and prevalence on the minimum  
 477 value. Importantly, very different interpretations of classification accuracy could be drawn from the  
 478 use of the kappa coefficient and overall accuracy. Even though 95% of the cases in the confusion  
 479 matrix have been correctly labelled it would be possible for a negative kappa coefficient to be  
 480 estimated that would indicate the level of agreement was less than that due to chance. While the  
 481 minimum kappa coefficient could be usefully interpreted as highlighting a poor classification, with  
 482 virtually all cases allocated to one class and the accuracy for one class zero, intermediate values could  
 483 be obtained. For example, Figure 6 shows one matrix for which the overall accuracy and producer's  
 484 accuracy for each class are all approximately 95%, highlighting a very accurate classification. The  
 485 kappa coefficient for the matrix in Figure 6 is 0.592 which lies in the range of 'moderate' agreement  
 486 in the Landis and Koch (1977) scale yet the classification meets an exacting Anderson-type target of  
 487 an overall accuracy of 95% with a producer's accuracy of at least 95% for each class; note purely for  
 488 ease of argument the focus is on the accuracy estimate itself relative to the target value and not its  
 489 associated confidence interval although the use of the latter may sometimes be appropriate.

490

491

475	50	525
0	475	475
475	525	1000

0	25	25
25	950	975
25	975	1000

492

493

(a)

(b)

494 Figure 5. Example confusion matrices to illustrate the range of possible kappa coefficients that could  
 495 arise for a classification with  $p_o=95\%$  (Figure 4). The layout of the matrices is as defined in Figure 1  
 496 and a sample of 1000 cases assumed. (a) Matrix for the maximum possible kappa coefficient,  $\kappa =$   
 497 0.900 (95% CI 0.873 - 0.927). (b) Matrix for the minimum possible kappa coefficient,  $\kappa = -0.026$   
 498 (95% CI -0.033 - -0.019).

499

40	47	87
2	911	913
42	958	1000

500

501 Figure 6. Confusion matrix for a classification that meets an Anderson-type target of an overall  
502 accuracy  $\geq 95\%$  and the producer's accuracy for each class are approximately equal and  $\geq 95\%$ . For  
503 this matrix,  $p_o = 95.1\%$ , and the producer's accuracies are 95.23% and 95.09%. The kappa coefficient  
504 for this matrix is  $\kappa = 0.592$  (95% CI 0.496 – 0.698).

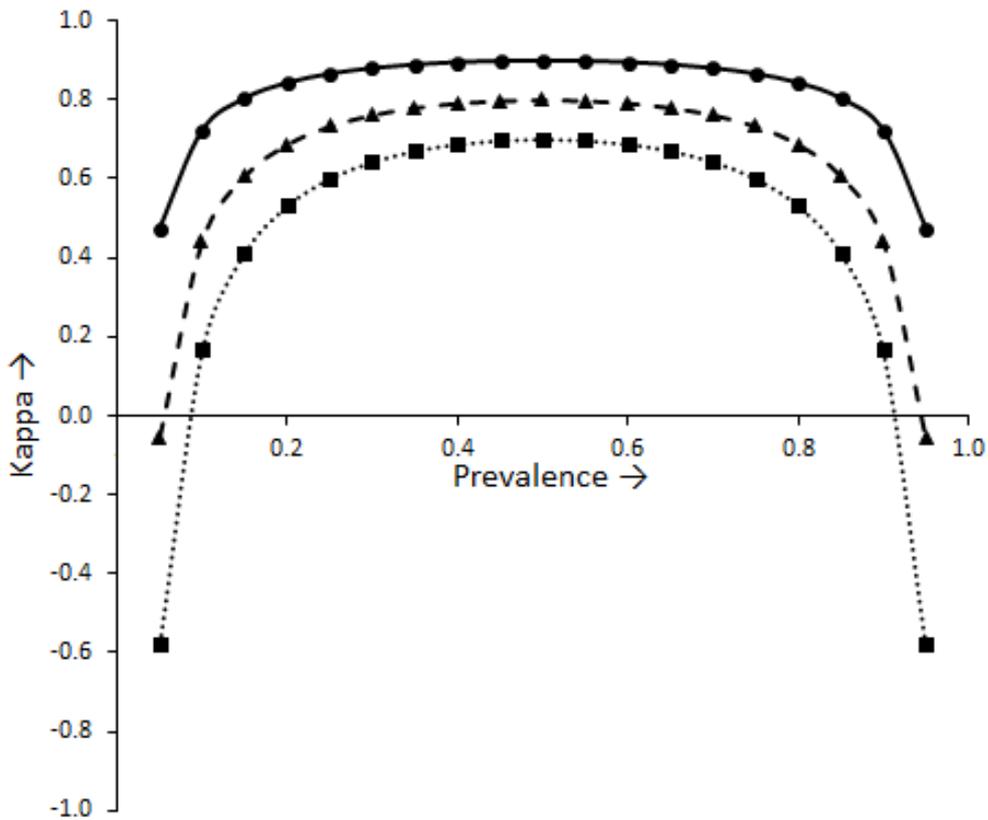
505

506

507 A key concern with the use of the kappa coefficient is its prevalence dependency (Byrt et al., 1993;  
508 Feinstein and Cicchetti, 1990; Sim and Wright, 2005). Again, while this is well-known it may be that  
509 the size of the effect is not fully appreciated. Figure 7 shows how the magnitude of the kappa  
510 coefficient varies with prevalence for three scenarios with a fixed overall accuracy (Vach, 2005):  
511 overall accuracies of 85%, 90% and 95%. Note the magnitude of the kappa coefficient varies greatly  
512 and the effects of prevalence are especially apparent at very large or low values of prevalence. In  
513 addition, a single value for the kappa coefficient could be associated with classifications of different  
514 overall accuracy due to differences in prevalence. Indeed differences in prevalence could change the  
515 apparent order or ranking of a series of classifications. For example, a classification could be viewed  
516 as being more accurate than another in terms of overall accuracy yet the exact opposite trend could be  
517 provided by the kappa coefficients; ranking classifications in terms of accuracy requires careful  
518 interpretation. The effect of prevalence variations is also very large and is further illustrated in Figure  
519 8 which shows matrices for four scenarios in which the overall accuracy and producer's accuracy for  
520 each class are fixed at 90% but which differ in prevalence. Each of the four matrices shown in Figure  
521 8 have the same overall accuracy and producer's accuracies but the magnitude of the kappa  
522 coefficient differs greatly. Indeed the 95% confidence intervals fitted to the four estimates of the  
523 kappa coefficient only just touch for two of the scenarios shown (Figure 8b and 8c). Comparing kappa  
524 coefficients is, therefore, challenging if there are differences in prevalence. Thus, the kappa

525 coefficient would not be a suitable measure if comparing classifications of study areas that may  
526 contain the same classes but at different abundances; similar problems with prevalence dependency  
527 may be observed with many other measures of accuracy. Would a difference in the magnitude of  
528 observed kappa coefficients indicate a difference in the quality of class labelling or merely reflect the  
529 variations in class prevalence?

530



531

532

533 Figure 7. Variation in the magnitude of the kappa coefficient with prevalence for three fixed value of  
534 overall accuracy. Three scenarios are shown in which the marginal values (i.e.  $n_{+}$  and  $n_{+}$ ) are equal  
535 and the overall accuracy is 85% (dotted line with square symbols), 90% (dashed line with triangular  
536 symbols) and 95% (solid line with circular symbols).

537

450	50	500	Prevalence = 0.50 $\kappa = 0.800$ (95% CI 0.763 - 0.837)
50	450	500	
500	500	1000	

(a)

90	90	180	Prevalence = 0.10 $\kappa = 0.590$ (95% CI 0.520 - 0.661)
10	810	820	
100	900	1000	

(b)

45	95	140	Prevalence = 0.05 $\kappa = 0.432$ (95% CI 0.344 - 0.520)
5	855	860	
50	950	1000	

(c)

9	99	108	Prevalence = 0.01 $\kappa = 0.137$ (95% CI 0.055 - 0.218)
1	891	892	
10	990	1000	

(d)

538

539 Figure 8. Confusion matrices for a scenario in which there is constant agreement on an overall and  
 540 per-class basis ( $p_o = 0.9$ , producer's accuracy for each class = 90%) but varying prevalence. (a)  
 541 prevalence = 0.5 (i.e. the two classes have equal abundance), (b) prevalence = 0.10, (c) prevalence =  
 542 0.05, and (d) prevalence = 0.01.

543

544 The various problems associated with the interpretation of the kappa coefficient make comparison of  
 545 kappa coefficients difficult, especially if the comparison is between studies of regions of dissimilar  
 546 prevalence (Usherax, 1987; Byrt et al., 1993; Vach, 2005; Sim and Wright, 2005). A major concern is  
 547 that the magnitude of a kappa coefficient and its possible range of values reflect the nature of the  
 548 population being studied (e.g. prevalence) (Byrt et al., 1993; Lantz and Nebenzahl, 1996). The kappa  
 549 coefficient has been widely promoted as a summary statistic that is meant to convey information on  
 550 thematic accuracy but it is a poor tool as it is highly mis-leading (Maclure and Willett, 1987). The  
 551 kappa coefficient is not well suited for use in accuracy assessment. Rather than use the kappa  
 552 coefficient because other studies have done so, and perpetuate a mistake, researchers should select an  
 553 accuracy measure appropriate for the task in-hand recognising that different measures of accuracy

554 reflect different aspects of quality and may require careful interpretation. Inspired by the comments of  
555 the referees on this article, as part of an effective peer review process, referees and editors should  
556 perhaps challenge the use of a measure such as the kappa coefficient in applications such as accuracy  
557 assessment and comparison for which it is unsuitable.

558

559 Finally on the issue of prevalence, it may be worth remembering that at the outset one key reason for  
560 not using overall accuracy was because of its sensitivity to the effect of variations in prevalence. This  
561 dependency is well known with  $p_o = (\theta P_+ + (1-\theta)P_-)$ . Overall accuracy is certainly an imperfect  
562 measure, as is any omnibus index (Byrt et al., 1993; Cicchetti and Feinstein, 1990), and no single  
563 measure will be universally ideal for accuracy assessment (Stehman, 1997a) but the kappa coefficient  
564 does not solve the problems associated with overall accuracy. That the kappa coefficient is prevalent  
565 dependent should come as no surprise given it is calculation from  $p_o$  and  $p_e$  in equation 1. Kappa is  
566 simply a rescaled version of  $p_o$  and  $p_e$  is prevalent dependent as prevalence is included in its  
567 calculation (equation 3). Because of the limitations of overall accuracy researchers have been  
568 encouraged to state per-class accuracies, such as user's and producer's accuracy, in addition (e.g. Liu  
569 et al., 2007; Stehman, 2000; Olofsson et al., 2014). A further enhancement would be to follow further  
570 good practices such as the provision of the confusion matrix and details of the sample used in its  
571 construction to allow estimation of other measures, even the kappa coefficient, if desired (Olofsson et  
572 al., 2013, 2014). It is difficult to identify how the provision of the kappa coefficient adds positively to  
573 this situation. The kappa coefficient alone is mis-leading so other information, notably on bias and  
574 prevalence, needs to be provided with it. The provision of a difficult to interpret measure such as the  
575 kappa coefficient that must be accompanied by additional measures such as bias and prevalence to aid  
576 interpretation does not help communicate accuracy information in a clear and succinct way. Then, in  
577 addition, there are concerns about the way chance is modelled and used. Given that the kappa  
578 coefficient is estimated from overall accuracy, it is evident that the estimation of the kappa coefficient  
579 is an unhelpful and unnecessary step in the assessment or comparison of classification accuracy.

580

#### 581 4. Conclusions

582 The kappa coefficient is widely promoted and used as a measure of thematic accuracy in remote  
583 sensing. The publications that promoted the use of the kappa coefficient have played an enormously  
584 influential role to inspire thought concerning rigorous quantitative assessments of classifications but  
585 promoted an inappropriate index. The reasons espoused for the use of the kappa coefficient are flawed  
586 and/or irrelevant as they apply equally well to other measures. Critically, the kappa coefficient is not  
587 an index of accuracy but a measure of the level of agreement observed beyond chance that is obtained  
588 using a model of chance that is inappropriate to the typical accuracy assessment scenario. Not only is  
589 the effect of chance agreement mis-estimated it is, however, irrelevant to an accuracy assessment  
590 which seeks to indicate the amount of error, and thereby correctness, in the labelling with the source  
591 of error inconsequential. The kappa coefficient is an inappropriate index to use to describe  
592 classification accuracy.

593

594 Many of the concerns with the kappa coefficient have been known for decades and it may be that its  
595 continued use in remote sensing is, in part, because the problems are viewed as being small and  
596 insubstantial. Here, emphasis has been placed on indicating the size and nature of the problems with  
597 the kappa coefficient by showing how its magnitude can vary as a function of basic properties of a  
598 study such as prevalence. Critically, simple examples have been used to show the unsuitability of the  
599 kappa coefficient for the description of accuracy and its comparison. For example, it was shown that  
600 classifications with an overall accuracy of 95% could have a kappa coefficient that lay within the  
601 range from -0.026 to 0.900. The difficulty of interpreting the estimated kappa coefficients is further  
602 highlighted by noting that the entire spread of possible values covers the complete range of the widely  
603 used Landis and Koch (1977) interpretation scale. Furthermore, if the classification satisfied a  
604 demanding Anderson-type target that required the producer's accuracy for each class be  $\geq 95\%$  the  
605 kappa coefficient for this very accurate classification would be interpreted as showing only moderate  
606 agreement. A key problem is the effect of variations in class abundance or prevalence, the very  
607 problem highlighted in criticisms of overall accuracy. Differences in prevalence make the comparison

608 of kappa coefficients very difficult, a researcher will be unsure if a difference reflects dissimilarity in  
609 the level of agreement or of the populations being studied. Overall accuracy on the other hand, while  
610 flawed, does have a clear meaning and, relative to kappa, is simple to estimate.

611

612 Different measures of accuracy reflect different aspects of a classification (Hand, 2012). Care must,  
613 therefore, be taken to ensure that a measure of accuracy that is appropriate for the task in-hand is  
614 adopted. There are many possible motivations and interests in an accuracy assessment which makes  
615 the provision of universal recommendations difficult. The literature on accuracy assessment can at  
616 times be challenging and other researchers may be better qualified to comment with authority and  
617 clarity on the topic but the common practice of using the kappa coefficient to indicate classification  
618 accuracy is flawed. Indeed, from the discussion above it is recommended that the kappa coefficient be  
619 dropped from the community's toolbox or at least used only sparingly and when good reason for its  
620 estimation exists such as in the assessment of agreement in class labelling among multiple  
621 interpreters. Although there are sometimes challenges to fully documenting an accuracy assessment,  
622 the provision of overall accuracy and per-class accuracy values together with the confusion matrix, set  
623 in the context of broader good practices (e.g. Olofsson et al., 2014; Stehman and Foody, 2019), should  
624 meet the objectives of most accuracy assessments. The provision of such information also allows  
625 assessments from other perspectives and the estimation of other measures, including even the kappa  
626 coefficient if desired, in order to meet the specific aims of a study. Comparisons of accuracy  
627 statements can be undertaken using overall accuracy and per-class accuracy using the same approach  
628 suggested for kappa if the samples involved are independent. If the samples are not independent, as is  
629 often the case in remote sensing research, alternative means to compare classification accuracy such  
630 as the McNemar test may be used. The kappa coefficient does not add positively to such accuracy  
631 assessments and comparisons. Given the challenges with its interpretation, the kappa coefficient  
632 should, therefore, not be used and reported routinely.

633

634 **Acknowledgments**

635 This article draws on experiences in accuracy assessment over many years that have benefited from  
636 inputs by colleagues ranging from article referees through to research collaborators which is much  
637 appreciated. The highly constructive and thoughtful comments from the three referees and editors are  
638 greatly appreciated and helped enhance the final version of the article.

639

640 **References**

641 Allouche, O., Tsoar, A. and Kadmon, R., 2006. Assessing the accuracy of species distribution models:  
642 prevalence, kappa and the true skill statistic (TSS). *Journal of Applied Ecology*, 43(6), 1223-1232.

643

644 Anderson, J. R., 1971. Land-use classification schemes, *Photogrammetric Engineering*, 37,  
645 379-387.

646

647 Anderson, J. R., Hardy, E. E., Roach, J. T. and Witmer, R. E., 1976. *A Land Use  
648 and Land Cover Classification System for Use with Remote Sensor Data*, Geological Survey  
649 Professional Paper 964, 28pp.

650

651 Ariza-López, F.J., Rodríguez-Avi, J., Alba-Fernández, M.V. and García-Balboa, J.L., 2019. Thematic  
652 accuracy quality control by means of a set of multinomials. *Applied Sciences*, 9, 4240.

653

654 Banerjee, M., Capozzoli, M., McSweeney, L. and Sinha, D., 1999. Beyond kappa: A review of  
655 interrater agreement measures. *Canadian Journal of Statistics*, 27(1), 3-23.

656

657 Brennan, R. L., and Prediger, D. J., 1981. Coefficient kappa: some uses, misuses, and alternatives.  
658 *Educational and Psychological Measurement*, 41, 687-699.

659

660 Byrt, T., Bishop, J. and Carlin, J.B., 1993. Bias, prevalence and kappa. *Journal of Clinical*  
661 *Epidemiology*, 46(5), 423-429.

662

663 Cicchetti, D.V. and Feinstein, A.R., 1990. High agreement but low kappa: II. Resolving the  
664 paradoxes. *Journal of Clinical Epidemiology*, 43(6), 551-558.

665

666 Cohen, J., 1960. A coefficient of agreement for nominal scales. *Educational and Psychological*  
667 *Measurement*, 20(1), 37-46.

668

669 Comber, A., Fisher, P., Brunson, C. and Khmag, A., 2012. Spatial analysis of remote sensing image  
670 classification accuracy. *Remote Sensing of Environment*, 127, 237-246.

671

672 Congalton, R.G., 1991. A review of assessing the accuracy of classifications of remotely sensed data.  
673 *Remote Sensing of Environment*, 37(1), 35-46.

674

675 Congalton, R.G. and Green, K., 2009. *Assessing the Accuracy of Remotely Sensed Data: Principles*  
676 *and Practices*. Second edition, CRC press, Boca Raton.

677

678 Congalton, R.G. and Mead, R.A., 1983. A quantitative method to test for consistency and correctness  
679 in photointerpretation. *Photogrammetric Engineering and Remote Sensing*, 49(1), 69-74.

680

681 Congalton, R.G., Oderwald, R.G. and Mead, R.A., 1983. Assessing Landsat classification accuracy  
682 using discrete multivariate analysis statistical techniques. *Photogrammetric Engineering and Remote*  
683 *Sensing*, 49(12), 1671-1678.

684

685 Czaplewski, R. L. 1994. *Variance Approximations for Assessments of Classification Accuracy*. Res.  
686 Pap. RM-316. Fort Collins, CO: U.S. Department of Agriculture, Forest Service, Rocky Mountain  
687 Forest and Range Experiment Station. 29 p.

688

689 Donner, A., Shoukri, M.M., Klar, N. and Bartfay, E., 2000. Testing the equality of two dependent  
690 kappa statistics. *Statistics in Medicine*, 19(3), 373-387.

691

692 Feinstein, A.R. and Cicchetti, D.V., 1990. High agreement but low kappa: I. The problems of two  
693 paradoxes. *Journal of Clinical Epidemiology*, 43(6), 543-549.

694

695 Fielding, A.H. and Bell, J.F., 1997. A review of methods for the assessment of prediction errors in  
696 conservation presence/absence models. *Environmental Conservation*, 24(1), 38-49.

697

698 Finn, J.T., 1993. Use of the average mutual information index in evaluating classification error and  
699 consistency. *International Journal of Geographical Information Science*, 7(4), 349-366.

700

701 Fitzgerald, R.W. and Lees, B.G., 1994. Assessing the classification accuracy of multisource remote  
702 sensing data. *Remote Sensing of Environment*, 47(3), 362-368.

703

704 Fleiss, J.L., Cohen, J. and Everitt, B.S., 1969. Large sample standard errors of kappa and weighted  
705 kappa. *Psychological Bulletin*, 72(5), p.323.

706

707 Fleiss, J.L., Levin, B. and Paik, M.C., 2013. *Statistical Methods for Rates and Proportions*. Third  
708 edition, John Wiley & Sons.

709

710 Foody, G.M., 1992. On the compensation for chance agreement in image classification accuracy  
711 assessment, *Photogrammetric Engineering and Remote Sensing*, 58, 1459-1460.

712

713 Foody, G.M., 2004. Thematic map comparison: evaluating the statistical significance of differences in  
714 classification accuracy. *Photogrammetric Engineering and Remote Sensing*, 70(5), 627-633.

715

716 Foody, G.M., 2008. Harshness in image classification accuracy assessment. *International Journal of*  
717 *Remote Sensing*, 29(11), 3137-3158.

718

719 Foody, G.M., 2009. Classification accuracy comparison: hypothesis tests and the use of confidence  
720 intervals in evaluations of difference, equivalence and non-inferiority. *Remote Sensing of*  
721 *Environment*, 113(8), 1658-1663.

722

723 Foody, G.M., 2010. Assessing the accuracy of land cover change with imperfect ground reference  
724 data. *Remote Sensing of Environment*, 114(10), 2271-2285.

725

726 Foody, G.M., 2011. Latent class modeling for site-and non-site-specific classification accuracy  
727 assessment without ground data. *IEEE Transactions on Geoscience and Remote Sensing*, 50(7), 2827-  
728 2838.

729

730 Hand, D. J., 2012. Assessing the performance of classification methods, *International Statistical*  
731 *Review*, 80 (3), 400-414.

732

733 Hoehler, F.K., 2000. Bias and prevalence effects on kappa viewed in terms of sensitivity and  
734 specificity. *Journal of Clinical Epidemiology*, 53(5), 499-503.

735

736 Hudson, W. D. and Ramm, C. W., 1987. Correct formulation of the Kappa coefficient of agreement,  
737 *Photogrammetric Engineering and Remote Sensing*, 53, 421-422.

738

739 Jansen, L L F, and van der Wel, F.J., 1994. Accuracy assessment of satellite derived landcover data:  
740 A review. *Photogrammetric Engineering and Remote Sensing*, 60(4), 479-426.

741

742 Jiang, S. and Liu, D., 2011. On chance-adjusted measures for accuracy assessment in remote sensing  
743 image classification. In *ASPRS Annual Conference, ASPRS 2011 Annual Conference Milwaukee,*  
744 *Wisconsin, May 1-5, 2011*

745

746 Landis, J.R. and Koch, G.G., 1977. The measurement of observer agreement for categorical data.  
747 *Biometrics*, 33, 159-174.

748

749 Lantz, C.A. and Nebenzahl, E., 1996. Behavior and interpretation of the  $\kappa$  statistic: Resolution of the  
750 two paradoxes. *Journal of Clinical Epidemiology*, 49(4), 431-434.

751

752 Liu, C., Frazier, P. and Kumar, L., 2007. Comparative assessment of the measures of thematic  
753 classification accuracy. *Remote Sensing of Environment*, 107(4), 606-616.

754

755 Maclure, M. and Willett, W.C., 1987. Misinterpretation and misuse of the kappa statistic. *American*  
756 *Journal of Epidemiology*, 126(2), 161-169.

757

758 Manel, S., Williams, H.C. and Ormerod, S.J., 2001. Evaluating presence–absence models in ecology:  
759 the need to account for prevalence. *Journal of Applied Ecology*, 38(5), 921-931.

760

761 Monserud, R.A. and Leemans, R., 1992. Comparing global vegetation maps with the Kappa statistic.  
762 *Ecological Modelling*, 62(4), 275-293.

763

764 Morales-Barquero, L., Lyons, M.B., Phinn, S.R. and Roelfsema, C.M., 2019. Trends in remote  
765 sensing accuracy assessment approaches in the context of natural resources. *Remote Sensing*, 11,  
766 2305.

767

768 Nishii, R. and Tanaka, S., 1999. Accuracy and inaccuracy assessments in land-cover classification.  
769 *IEEE Transactions on Geoscience and Remote Sensing*, 37(1), 491-498.

770

771 Olofsson, P., Foody, G.M., Stehman, S.V. and Woodcock, C.E., 2013. Making better use of accuracy  
772 data in land change studies: estimating accuracy and area and quantifying uncertainty using stratified  
773 estimation. *Remote Sensing of Environment*, 129, 122-131.

774

775 Olofsson, P., Foody, G.M., Herold, M., Stehman, S.V., Woodcock, C.E. and Wulder, M.A., 2014.  
776 Good practices for estimating area and assessing accuracy of land change. *Remote Sensing of*  
777 *Environment*, 148, 42-57.

778

779 Pontius Jr, R.G., 2000. Comparison of categorical maps. *Photogrammetric Engineering and Remote*  
780 *Sensing*, 66(8), 1011-1016.

781

782 Pontius Jr, R.G. and Millones, M., 2011. Death to Kappa: birth of quantity disagreement and  
783 allocation disagreement for accuracy assessment. *International Journal of Remote Sensing*, 32(15),  
784 4407-4429.

785

786 Pontius, R.G. and Parmentier, B., 2014. Recommendations for using the relative operating  
787 characteristic (ROC). *Landscape Ecology*, 29(3), 367-382.

788

789 Rogan, W. J. and Gladen, B. (1978) Estimating prevalence from the results of a screening  
790 test, *American Journal of Epidemiology*, 107, 71-76.

791

792 Rosenfield, G.H. and Fitzpatrick-Lins, K., 1986. A coefficient of agreement as a measure of thematic  
793 classification accuracy. *Photogrammetric Engineering and Remote Sensing*, 52(2), 223-227.

794 Sim, J. and Wright, C.C., 2005. The kappa statistic in reliability studies: use, interpretation, and  
795 sample size requirements. *Physical Therapy*, 85(3), 257-268.

796

797 Smits, P.C., Dellepiane, S.G. and Schowengerdt, R.A., 1999. Quality assessment of image  
798 classification algorithms for land-cover mapping: A review and a proposal for a cost-based approach.  
799 *International Journal of Remote Sensing*, 20(8), 1461-1486.

800

801 Stehman, S., 1996. Estimating the kappa coefficient and its variance under stratified random  
802 sampling. *Photogrammetric Engineering and Remote Sensing*, 62(4), 401-407.

803

804 Stehman, S.V., 1997a. Selecting and interpreting measures of thematic classification accuracy.  
805 *Remote sensing of Environment*, 62(1), 77-89.

806

807 Stehman, S.V., 1997b. Estimating standard errors of accuracy assessment statistics under cluster  
808 sampling. *Remote Sensing of Environment*, 60(3), 258-269.

809

810 Stehman, S.V., 2000. Practical implications of design-based sampling inference for thematic map  
811 accuracy assessment. *Remote Sensing of Environment*, 72(1), 35-45.

812

813 Stehman, S.V. and Foody, G.M., 2009. Accuracy assessment. In: Warner, T A, Nellis, N. D. and  
814 Foody, G M.(editors) *The SAGE Handbook of Remote Sensing*, Sage, London, 297-309.

815

816 Stehman, S.V. and Foody, G.M., 2019. Key issues in rigorous accuracy assessment of land cover  
817 products. *Remote Sensing of Environment*, 231, 111199.

818

819 Story, M. and Congalton, R.G., 1986. Accuracy assessment: a user's perspective. *Photogrammetric*  
820 *Engineering and Remote Sensing*, 52(3), 397-399.

821

822 Thron, C. and Miller, V., 2015. Persistent confusions about hypothesis testing in the social  
823 sciences. *Social Sciences*, 4(2), 361-372.

824 Tsutsumida, N. and Comber, A.J., 2015. Measures of spatio-temporal accuracy for time series land  
825 cover data. *International Journal of Applied Earth Observation and Geoinformation*, 41, 46-55.

826

827 Türk, G., 1979. GT index: A measure of the success of prediction. *Remote Sensing of Environment*,  
828 8(1), 65-75.

829

830 Türk, G. 2002. Map evaluation and 'chance correction'. *Photogrammetric Engineering and Remote*  
831 *Sensing*, 68, 123–133.

832

833 Uebersax, J.S., 1987. Diversity of decision-making models and the measurement of interrater  
834 agreement. *Psychological Bulletin*, 101(1), p.140.

835

836 Vach, W., 2005. The dependence of Cohen's kappa on the prevalence does not matter. *Journal of*  
837 *Clinical Epidemiology*, 58(7), 655-661.

838

839 Wu, S.M., Whiteside, U. and Neighbors, C., 2007. Differences in inter-rater reliability and accuracy  
840 for a treatment adherence scale. *Cognitive Behaviour Therapy*, 36(4), 230-239.

841

842 Ye, S., Pontius Jr, R.G. and Rakshit, R., 2018. A review of accuracy assessment for object-based  
843 image analysis: From per-pixel to per-polygon approaches. *ISPRS Journal of Photogrammetry and*  
844 *Remote Sensing*, 141, 137-147.