# Supervised methods of image segmentation accuracy assessment in land cover mapping

*Hugo Costa, Giles M. Foody, and Doreen S. Boyd*

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

## Abstract

Land cover mapping via image classification is sometimes realized through object-based image analysis. Objects are typically constructed by partitioning imagery into spatially contiguous groups of pixels through image segmentation and used as the basic spatial unit of analysis. As it is typically desirable to know the accuracy with which the objects have been delimited prior to undertaking the classification, numerous methods have been used for accuracy assessment. This paper reviews the state-of-the-art of image segmentation accuracy assessment in land cover mapping applications. First the literature published in three major remote sensing journals during 2014-2015 is reviewed to provide an overview of the field. This revealed that qualitative assessment based on visual interpretation was a widely-used method, but a range of quantitative approaches is available. In particular, the empirical discrepancy or supervised methods that use reference data for assessment are thoroughly reviewed as they were the most frequently used approach in the literature surveyed. Supervised methods are grouped into two main categories, geometric and non-geometric, and are translated here to a common notation which enables them to be coherently and unambiguously described. Some key considerations on method selection for land cover mapping applications are provided, and some research needs are discussed.

1

# 1 Introduction

Land cover mapping is a very common application of remote sensing and has been increasingly conducted through object-based image analysis (Blaschke, 2010). Object-based image analysis has been described as an advantageous alternative to conventional per-pixel image classification, and adopted in a diverse range of studies (Bradley, 2014; Feizizadeh et al., 2017; Matikainen et al., 2017; Strasser and Lang, 2015).

Objects are typically discrete and mutually exclusive groups of neighbouring pixels and used as the basic spatial unit of analysis. Objects may be delimited or obtained via a range of sources (e.g. cadastral data), but typically are constructed through an image segmentation analysis, and thus often called segments. In this paper the terms "object" and "segment" are used synonymously. Image segmentation is performed by algorithms with the purpose of constructing objects corresponding to geographical features distinguishable in the remotely sensed data, which may be useful for applications such as land cover mapping.

Constructing objects poses a set of challenges. For example, it is necessary to select a segmentation algorithm from the numerous options available, but comparative studies (e.g. Basaeed et al., 2016; Neubert et al., 2008) are uncommon. Also each of the segmentation algorithms is typically able to produce a vast number of outputs depending on the parameter settings used. Selecting the most appropriate segmentation is, therefore, difficult.

Multiple methods have been proposed to assess the accuracy of an image segmentation and are normally grouped in two main categories: empirical discrepancy and empirical goodness methods, also commonly referred to as supervised and unsupervised methods respectively (Zhang, 1996). Most of the supervised methods essentially compare a segmentation output to a reference data set and measure the similarity or discrepancy between the two

2

46  representations (e.g. overlapping area) (Clinton et al., 2010). Unsupervised methods measure

47  some desirable properties of the segmentation outputs (e.g. object's spectral homogeneity),

48  thus measuring their quality (Zhang et al., 2008).

49  There is no standard approach for image segmentation accuracy assessment, and some studies

50  have compared accuracy assessment methods. Supervised and unsupervised methods are

51  normally compared separately. For example, with regard to supervised methods, Clinton et al.

52  (2010), Räsänen et al. (2013), and Whiteside et al. (2014) compared dozens of methods, all of

53  them focused on some geometric property of the objects, such as positional accuracy relative

54  to the reference data. These and other studies highlight the differences and similarities

55  obtained from the methods compared so the reader gains a perspective of the field. However,

56  many other supervised methods have been proposed yet are barely compared against previous

57  counterparts; these tend to be newly proposed methods (e.g. Costa et al., 2015; Liu and Xia,

58  2010; Marpu et al., 2010; Su and Zhang, 2017). Furthermore, the methods are often described

59  using a notation suitable for the specific case under discussion, which makes the cross-

60  comparison of methods difficult.

61  Studies like Clinton et al. (2010) are valuable in reviewing the field of image segmentation

62  accuracy assessment, but they often focus on the geometry of the objects evaluated and

63  ignore that a supervised but non-geometric approach may be followed (e.g. Wang et al.

64  2004). Moreover, supervised methods are typically compared within a specific study case

65  without discussion of further and important issues, such as the suitability of the methods as a

66  function of context. As image segmentation is increasingly used in a wide range of

67  applications, the behaviour and utility of specific methods is expected to vary in each case.

68  Thus, selecting a method to assess the accuracy of image segmentation may be based on an

69  incomplete understanding of the available options and ultimately problematic.

70 This paper reviews the state-of-the-art of image segmentation accuracy assessment in land

71 cover mapping applications. The literature published in three major remote sensing journals

72 in 2014-2015 is reviewed to provide an overview of the field, namely the methods used and

73 their popularity. In particular, the supervised methods are thoroughly reviewed as they are

74 widely used. A comprehensive description of which supervised methods are available is

75 presented with the aim of providing a basis on which the remote sensing community may

76 consider and select a suitable method for particular applications. A discussion on which

77 methods should be used is provided, and research needs are highlighted.

78 **2  Background**

79 Image objects are typically expected to delimit features of the Earth's surface such as land

80 cover patches that are remotely sensed using an air/spaceborne imaging system. Image

81 segmentation cannot, however, deliver results exactly according to the desired outcome for

82 multiple reasons, such as unsuitable definition of segmentation algorithm parameter settings,

83 and insufficient spectral and spatial resolution of the data. Thus, image segmentation error is

84 common, namely under- and over-segmentation. Under-segmentation error occurs when

85 image segmentation fails to define individual objects to represent different contiguous land

86 cover classes, thus constructing a single object that may contain more than one land cover

87 class. On the contrary, over-segmentation error occurs when unnecessary boundaries are

88 delimited, and thus multiple contiguous objects, potentially of the same land cover class, are

89 formed.

90 Segmentation errors have been traditionally identified through visual inspection, but it has

91 some drawbacks, especially when assessing large areas and comparing numerous

92 segmentation outputs. Specifically, visual interpretation is time consuming, subjective, and

4

93    the results produced by the same or different operators may not be reproducible (Coillie et al.,

94    2014; Lang et al., 2010). As a result, objective and quantitative methods for the assessment of

95    image segmentation accuracy may be necessary and have become more popular in recent

96    years.

97    The literature published during 2014-2015 in three remote sensing journals was reviewed to

98    provide an overview of the state-of-the-art of image segmentation accuracy assessment. The

99    journals were *Remote Sensing of Environment*, *ISPRS Journal of Photogrammetry and*

100   *Remote Sensing*, and *Remote Sensing Letters*. These journals were selected to represent the

101   variety of current publication outlets in the field. Historically, the former journal has had the

102   greatest impact factor among the remote sensing journals. The second journal has been

103   particularly active in publishing papers on object-based image analysis. The latter journal is a

104   relatively young journal dedicated to rapid publications. The papers that included specific

105   terms (namely "obia", "geobia", "object-based", and "object-oriented") in the title, abstract,

106   and key words were retained for analysis. A total of 55 out of 67 papers that matched the

107   search terms were identified as relevant, each describing techniques for estimating objects

108   which were used as the basic spatial unit in land cover mapping applications.

109   These 55 papers were analysed, and it was noticeable that 17 papers (30.9%) do not

110   document if or how the accuracy of the image segmentation outputs was assessed. This

111   shows that image segmentation accuracy assessment is often overlooked as an important

112   component of an image segmentation analysis protocol. It is speculated that visual

113   interpretation was used in most of the cases that provide no information accuracy, as having

114   used no sophisticated method may reduce any motivation for documenting the topic. The

115   remaining 38 papers explicitly described the methods used, and often more than one method

116   was adopted. Visual interpretation was widely used, with 15 papers (25.3% of the total of

5

117  papers) describing that the qualitative appearance of the segmentations influenced the

118  assessment of the results (e.g. Qi et al., 2015). Details were typically not given, such as the

119  time dedicated to visual interpretation and number of interpreters.

120  When a quantitative alternative to subjective visual interpretation was explicitly adopted, the

121  methods used varied widely. A rudimentary strategy of assessing the accuracy of image

122  segmentations, and used in five papers (9.1%), was to use simple descriptive statistics, such

123  as the average of some attributes of the objects like area, to get an impression of the

124  segmentation output. The statistics were used in a supervised or unsupervised fashion. In the

125  former situation, the statistics were compared to the statistics of a reference data set depicting

126  desired polygonal shapes, and small differences were regarded as indicative of large

127  segmentation accuracy (e.g. Liu et al., 2015). When no reference data were used (i.e.

128  unsupervised fashion), the statistics identified the image segmentation from the set obtained

129  with the most desirable properties, such as a target mean size (i.e. area) of the objects

130  (Hultquist et al., 2014). Although descriptive statistics can measure some quantitative

131  properties of an image segmentation, they provide a very limited sense of the accuracy of the

132  objects, for example in the spatial domain, and here they are not regarded as a true accuracy

133  assessment method. The latter are typically more evolved and normally grouped into

134  supervised and unsupervised methods.

135  Supervised methods were found in 21 (38.2%) of the papers reviewed (e.g. Zhang et al.,

136  2014). Although there was no dominant method, the Area Fit Index (Lucieer and Stein, 2002)

137  and Euclidean distance 2 (Liu et al., 2012) were the supervised methods that were most used

138  with three appearances each (Belgiu and Drăguţ, 2014; Drăguţ et al., 2014; Witharana et al.,

139  2014; Witharana and Civco, 2014; Yang et al., 2014). Many of the other methods identified

140  were used only once (e.g. Carleer et al. 2005). These and other supervised methods are,

6

141  however, thoroughly described in the next section. Unsupervised methods were applied in 13

142  (23.6%) of the papers surveyed (e.g. Robson et al., 2015). The unsupervised method most

143  used in the literature reviewed was the Estimation of Scale Parameter (ESP or ESP2) tool

144  (Drăguţ et al., 2014, 2010) available in the popular eCognition software. The segmentation

145  algorithms available in this software were used in most of the papers surveyed (36 papers,

146  65.5%) to construct image objects.

147  Object-based image analysis has received much attention and acceptance (Blaschke et al.,

148  2014; Dronova, 2015), but the accuracy assessment of image segmentation, which is a central

149  stage of the analysis, appears to be in a relatively early stage of maturation. Although

150  procedures for image segmentation accuracy assessments have not been standardized, a more

151  harmonized approach is desirable. Using subjective visual interpretation may be acceptable

152  and suitable for some applications; the reasons are seldom explained in the literature. Among

153  the quantitative methods proposed for image segmentation accuracy assessment, supervised

154  approaches seem to be the most frequently adopted, hence reviewed hereafter.

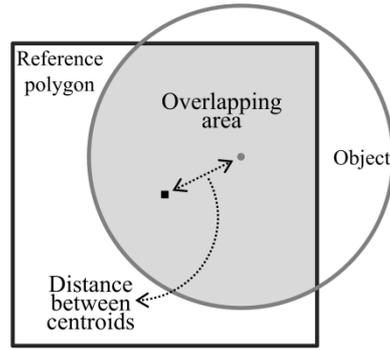155  **3  Supervised methods**

156  Supervised methods for image segmentation accuracy assessment use reference data to

157  estimate the accuracy of the objects constructed. Often the reference data are formed by

158  polygons extracted from the remotely sensed data in use (e.g. based on visual interpretation)

159  or collected externally (e.g. a field boundary map). Approaches for assessing accuracy based

160  on reference data are herein grouped into two main categories: geometric and non-geometric.

161  Geometric methods are the most widely used and typically focus on the geometry of the

162  objects and polygons to determine the level of similarity among them. Ideally, there should

163  be no difference among objects and polygons in terms of area, position, and shape. Note that

164    the land cover class(es) associated with the objects and polygons typically need not be

165    known.

166    With non-geometric methods the land cover class(es) associated with the objects must be

167    known, and reference data polygons are not always used. The properties of the objects such

168    as the spectral content are used in a variety of ways, depending on the specific method.

169    Ideally, the content of the objects representing different land cover classes should be as

170    different as possible. When polygons are also used, the content of objects and polygons

171    representing the same land cover class should be identical. Note that the spatial or geometric

172    correspondence between objects and polygons need not be known. Fuller details on both

173    geometric and non-geometric approaches are given in the sub-sections that follow.

174    Rudimentary strategies (for example used in 9.1% of the papers reviewed in the previous

175    section) are not covered however.

## 176   3.1   Geometric methods

177    Geometric methods rely on quantitative metrics that describe aspects of the geometric

178    correspondence between objects and polygons, often based on difference in area and position

179    (Winter, 2000). Figure 1 illustrates a typical case involving an object and polygon for which

180    the larger the overlapping area and/or the shorter the distance between their centroids, the

181    larger the accuracy with which the object has been delimited.

182
183 Figure 1. Geometric comparison between an object and polygon based on the overlapping

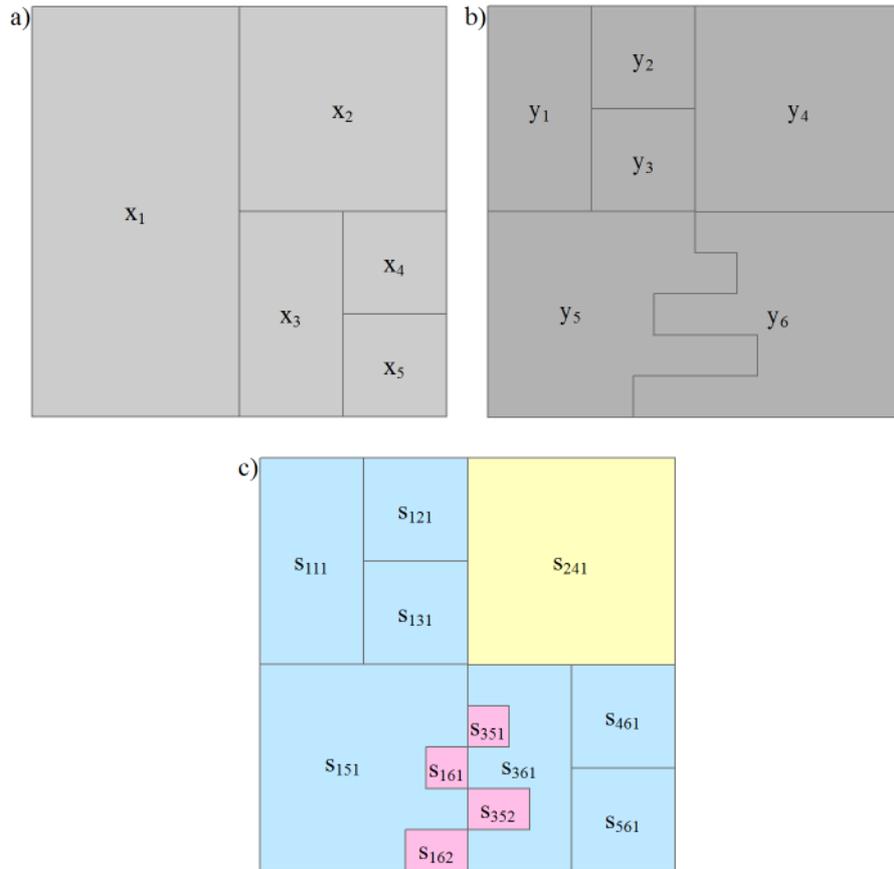184 area (shaded area) and/or distance between centroids (dashed arrow).

### 3.1.1 Notation

186 Notation is necessary to assist the description of the metrics used by geometric methods. The

187 notation presented hereafter uses that defined in Clinton et al. (2010). Therefore, the notation

188 is transcribed below together with additional elements necessary to describe all the methods

189 covered.

190 The $m$ objects constructed via image segmentation are denoted by $y_j$ ($j=1, \ldots, m$), the $n$

191 polygons forming a reference data set by $x_i$ ($i=1, \ldots, n$), and the $l$ pixels of the segmented

192 remotely sensed data by $z_p$ ($p=1, \ldots, l$). They define the following sets:

193 • $X=\{x_i: i=1, \ldots, n\}$ is the set of $n$ polygons (Figure 2a)
194 • $Y=\{y_j: j=1, \ldots, m\}$ is the set of $m$ objects of a segmentation output (Figure 2b)
195 • $S=X\cap Y=\{s_{ijk}: \text{area}(x_i\cap y_j)\neq0\}$ is the set of $s$ intersection objects that result from the
196   spatial intersection (represented by symbol $\cap$) of X and Y; $s_{ijk}$ is the $k^{th}$ object that
197   results from the spatial intersection of the $i^{th}$ polygon ($x_i$) with the $j^{th}$ object ($y_j$)
198   (Figure 2c)
199 • $Z=\{z_p: p=1,\ldots,l\}$ is the set of $l$ pixels of the segmented remotely sensed data.
200

201 Set S is the result of a spatial intersection of X and Y, which can be defined using common

202 geographical information systems. Note that the subscript k is needed to create a unique

203 symbol as the overlay of $x_i$ and $y_i$ can yield more than one discontinuous polygonal area ($x_1$

204 and $y_6$ in Figure 2). Set Z is simply the set of pixels that form the remotely sensed data

205 submitted to segmentation analysis, but its definition is nevertheless useful for describing

206 clearly some metrics.

9

207

Figure 2. Sets X, Y, and S: (a) reference set X, (b) segmentation Y, and (c) intersection S=X∩Y. In (c) yellow denotes one-to-one, blue denotes one-to-many, and pink denotes many-to-many (Section 3.1.1.3).

The description of the methods also requires the use of symbols that characterize the sets X, Y, S, and Z, and their members. For example, size() denotes the number of an item identified in brackets, for example the number of objects that belong to Y – size(Y) – or the number of pixels of an object – size($y_j$); and dist() is the distance between two items identified in brackets, for example the centroids of $y_j$ and $x_i$ – dist(centroid($x_i$), centroid($y_j$)). This basic notation is used to express more complex cases. For example, area($x_i \cap y_j$) is the area of the geographical intersection of polygon $x_i$ and object $y_j$. Other self-explanatory cases are used in the notation adopted. Furthermore, mathematical symbols are also used, such as ¬ which is the logical negation symbol and read as "not", \ which is the complement symbol used in set theory and reads as "minus" or "without", and ∪ which is the union symbol.

10

221 Subsets of X, Y, and S must be defined to assist the description of methods that follow four

222 different strategies: (i) Y is compared to X, (ii) X is compared to Y, (iii) S is compared to

223 both X and Y, and (iv) X and Y are compared to Z. In all of the cases, the definition of

224 subsets of X, Y, and S are used to decide which polygons $x_i$, objects $y_j$, and intersection

225 objects $s_{ijk}$ corresponds to each other or to pixel $z_p$, which is central to the calculation of

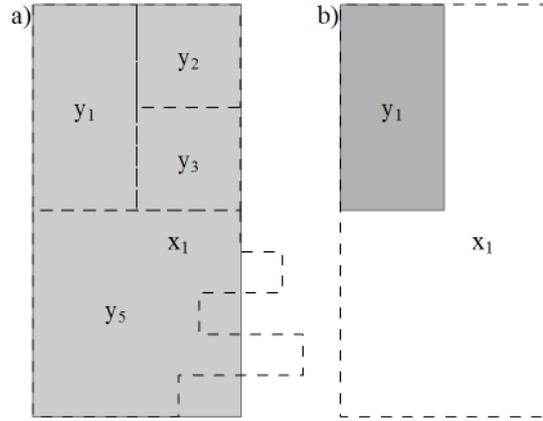226 geometric metrics (presented in Section 3.1.2).

227 **3.1.1.1 Set Y compared to set X**

228 In image segmentation accuracy assessment most often the set Y is compared to set X. This

229 strategy typically involves the calculation of geometric metrics for the members of X, and

230 thus there is the need to identify which member(s) of Y correspond to each member of X. For

231 example, Figure 3a shows the set of objects that overlap and thus can be considered as

232 corresponding to a polygon $x_i$. The specific objects that are actually considered as

233 corresponding depends on the method used, and the calculations related to each polygon $x_i$

234 consider only the objects regarded as corresponding. Thus, it is useful to define the following

235 subsets of Y for each member of X:

236 • $\tilde{Y}_i$ is the subset of Y such that $\tilde{Y}_i = \{y_j: \text{area}(x_i \cap y_j) \neq 0\}$

237 • $Ya_i$ is a subset of $\tilde{Y}_i$ such that $Ya_i = \{y_j: \text{the centroid of } x_i \text{ is in } y_j\}$

238 • $Yb_i$ is a subset of $\tilde{Y}_i$ such that $Yb_i = \{y_j: \text{the centroid of } y_j \text{ is in } x_i\}$

239 • $Yc_i$ is a subset of $\tilde{Y}_i$ such that $Yc_i = \{y_j: \text{area}(x_i \cap y_j)/\text{area}(y_j) > 0.5\}$

240 • $Yd_i$ is a subset of $\tilde{Y}_i$ such that $Yd_i = \{y_j: \text{area}(x_i \cap y_j)/\text{area}(x_i) > 0.5\}$

241 • $Ye_i$ is a subset of $\tilde{Y}_i$ such that $Ye_i = \{y_j: \text{area}(x_i \cap y_j)/\text{area}(y_j) = 1\}$

242 • $Yf_i$ is a subset of $\tilde{Y}_i$ such that $Yf_i = \{y_j: \text{area}(x_i \cap y_j)/\text{area}(y_j) > 0.55\}$

243 • $Yg_i$ is a subset of $\tilde{Y}_i$ such that $Yg_i = \{y_j: \text{area}(x_i \cap y_j)/\text{area}(y_j) > 0.75\}$

244 • $Y_i^* = Ya_i \cup Yb_i \cup Yc_i \cup Yd_i$

245 • $Y_i'$ is a subset of $\tilde{Y}_i$ such that $Y_i' = \{y_j: \max(\text{area}(x_i \cap y_j))\}$.

246 The definition of subsets of Y expresses the variety of criteria of correspondence that has

247 been used. For example, some methods require the centroid of the objects to fall inside the

248 polygons, and $Yb_i$ denotes the set of objects whose centroid falls inside a specific polygon $x_i$.

249 However, most of the criteria of correspondence used define a threshold of overlapping area

250 between polygons and objects. For example, at least half of the object's area may have to

251 overlap a polygon for a positive correspondence to be considered; $Yc_i$ denotes the set of

11

252 objects that comply with this criterion for a specific polygon $x_i$. The selection of a specific

253 subset of Y depends on the method used.



254
255 Figure 3. Comparison between X and Y of Figure 2: (a) four potential objects (dashed lines)

256 corresponding to polygon $x_1$ (grey background) when Y compares to X; (b) one potential

257 polygon (dashed line) corresponding to object $y_1$ (grey background) when X compares to Y.

**3.1.1.2 Set X compared to set Y**

259 When set X is compared to Y, geometric metrics are calculated for the members of Y, and

260 thus there is the need to identify which member(s) of X correspond to each member of Y. For

261 example, Figure 3b shows that one polygon overlap and thus can be considered as

262 corresponding to an object $y_j$. The calculations related to each object $y_j$ consider only the

263 polygons regarded as corresponding, depending on the method used. Thus, it is useful to

264 define the following subsets of X for each member of Y:

265 • $\tilde{X}_j$ is the subset of X such that $\tilde{X}_j = \{x_i: \text{area}(y_j \cap x_i) \neq 0\}$
266 • $Xc_j$ is a subset of $\tilde{X}_j$ such that $Xc_j = \{x_i: \text{area}(y_j \cap x_i)/\text{area}(y_j) > 0.5\}$
267 • $X'_j$ is a subset of $\tilde{X}_j$ such that $X'_j = \{x_i: \max(\text{area}(y_j \cap x_i))\}$
268 • $X''_j$ is a subset of $\tilde{X}_j$ such that $X''_j = \{x_i: \max(\text{area}(y_j \cap x_i)/\text{area}(y_j \cup x_i))\}$.
269

270 The subsets of X defined above represent the criteria of correspondence that have been used

271 when X is compared to Y. All the criteria define a threshold of overlapping area between

272 polygons and objects. For example, a polygon may have to overlap more than half of the

273 object's area for a positive correspondence between objects and polygons; $Xc_i$ denotes the set

12

274  of polygons that comply with this criterion for a specific object $y_j$. The selection of a specific

275  subset of X depends on the method used.

276  To describe two particular methods found in the literature (Costa et al., 2015; Liu and Xia,

277  2010), it is useful to define X not as the set of *n* reference polygons, but the set of *t* thematic

278  classes represented in X. For example, if $x_3$ and $x_4$ in Figure 2 are two polygons representing

279  the same thematic class, $c_i$, and both intersect the same object, $y_6$, notation like $area(c_i \cap y_6)$

280  can be used, where $area(c_i)=area(x_3 \cup x_4)$. Thus, similarly to above:

281  • C={$c_i$: i=1, …, *t*} is the set of *t* thematic classes represented in X; classes $c_i$ can also

282  be denoted as $d_i$ as it is useful to describe a specific method (Costa et al., 2015).

283  When comparing C to Y, the following subset of C is identified for each $y_j$:

284  • $\tilde{C}_j$ is the subset of C such that $\tilde{C}_j$ ={$c_i$: $area(c_i \cap y_j) \neq 0$}.

**3.1.1.3 Set S compared to both sets X and Y**

286  When set S is compared to both sets X and Y, three types of hierarchical relations between

287  polygons and objects emerge. The three types are one-to-one, one-to-many, and many-to-

288  many relations (Figure 2). The first type occurs when $x_i$ and $y_j$ match perfectly. One-to-many

289  relations occur when $x_i$ intersects several objects or *vice-versa*. Many-to-many relations occur

290  when several discontinuous intersection objects correspond to a same $x_i$ and $y_j$ (e.g. sliver

291  intersection objects $s_{ijk}$ along the edges of $x_i$ and $y_j$).

292  Given the three types of hierarchical object relations, the following subsets of S are defined:

293  • $S_1$={$s_{ijk}$: $area(x_i \cap y_j) = area(x_i \cup y_j)$} is the subset of all one-to-one objects
294  • $S_{2a}$={$s_{ijk}$: (one $x_i \cap$ many $y_j$) $\vee$ (many $x_i \cap$ one $y_j$)} is the subset of all one-to-many
295  relations
296  • $S_{2b}$={$s_{ijk}$: (one $x_i \cap$ many $y_j$) $\vee$ (many $x_i \cap$ one $y_j$); $max(area(s_{ijk}))$} is the subset of the
297  largest one-to-many relations
298  • $S_3$={$s_{ijk}$: one $x_i \cap$ one $y_j$ over discontinuous areas; $max(area(s_{ijk}))$} is the subset of the
299  largest many-to-many relations.

300  Based on the above subsets, it is useful to define the subsets $Sa=S_1 \cup S_{2a} \cup S_3$, and $Sb=S_1 \cup$

301  $S_{2b} \cup S_3$. Finally, subsets of Sa and Sb are defined for each $x_i$ and $y_j$:

302    • $Sax_i=\{s_{ijk}: area(s_{ijk} \cap x_i) \neq 0\}$
303    • $Say_j=\{s_{ijk}: area(s_{ijk} \cap y_j) \neq 0\}$
304    • $Sbx_i=\{s_{ijk}: area(s_{ijk} \cap x_i) \neq 0\}$
305    • $Sby_j=\{s_{ijk}: area(s_{ijk} \cap y_j) \neq 0\}$

306    The definition of subsets $Sax_i$, $Say_j$, $Sbx_i$, and $Sby_j$ are used in Möller et al. (2013) and Costa

307    et al. (2015).

### 3.1.1.4 Sets X and Y compared to set Z

309    To describe two particular methods found in the literature (Martin, 2003; Zhang et al., 2015a,

310    2015b), it is useful to consider the assessment framework at the pixel level and thus define

311    the following subsets of X and Y that correspond to each member of Z:

312    • $Xa_p$ is the subset of X such that $Xa_p=\{x_i:$ the centroid of $z_p$ is in $x_i\}$
313    • $Ya_p$ is the subset of Y such that $Ya_p=\{y_j:$ the centroid of $z_p$ is in $y_j\}$.

### 3.1.2  Available metrics

315    Geometric metrics are presented in Table 1 using the notation defined above, except four

316    cases that would require the definition of unnecessarily complex notation, and thus are

317    described as text (metrics 6, 7, 13 and 28). The metrics express the fundamental calculation

318    involving objects and polygons; each object, polygon, or intersection object receives a metric

319    value, which will tell something about the individual geometric accuracy of the objects

320    constructed. Assessing each areal entity individually is often referred to as local evaluation or

321    validation (Möller et al., 2013, 2007; Persello and Bruzzone, 2010). The subscripts i and j

322    used in the name of the metrics in Table 1 (e.g. $Precision_{ij}$) indicate that the metrics are

323    calculated for the local level. These subscripts come from those used to identify the specific

324    polygon $x_i$ and object $y_j$ involved in the calculations.

325    Place table 1 near here. **See Table 1 after the references.**

326    Local metric values are commonly aggregated in a variety of ways to produce a single value

327    to express the accuracy of a segmentation output as a whole. This is often referred to as

328    global evaluation or validation (Möller et al., 2013, 2007; Persello and Bruzzone, 2010).

329    Table 1 provides details on how the local metric values are aggregated for the global level in

330    the column headed Notes. Typically, the local values are summed or averaged in either one or

331    two steps, which in Clinton et al. (2010) is referred to as weighted and unweighted measures

332    respectively. In the first case, all the local values are aggregated in a straightforward fashion

333    (e.g. SimSize, metric 15). In the second case, the aggregation is undertaken first for each

334    individual polygon or object (depending of the strategy of comparison), and then for the

335    whole segmentation. For example, metric $PI_{ij}$ (metric 22) is first aggregated for each polygon,

336    and then for the whole segmentation. Therefore, if for a given polygon, say $x_1$, there are two

337    corresponding objects, $y_1$ and $y_2$, then $PI_{11}$ and $PI_{12}$ are calculated according to metric 22.

338    Then, $PI_{11}$ and $PI_{12}$ are summed to calculate a single $PI_1$ value for polygon $x_1$. This produces

339    $n$ $PI_i$ values (one for each polygon $x_i$). Finally, the $n$ $PI_i$ values can be averaged to express

340    image segmentation accuracy as a whole, denoted as PI (without any subscript).

341    Showing the metrics for the local level facilitates comparison, but it was not possible to write

342    them all in the same style. For example, the $LP_i$ formula (metric 31) shows only the subscript

343    i (i.e. the subscript j is missing). This specific metric, calculated for polygons $x_i$, needs

344    immediately to involve all the corresponding objects. In other cases, such as NSR (metric

345    39), the metric's name in Table 1 shows no subscripts because the metric is calculated

346    directly as a global value for the whole segmentation output.

347    Oftentimes the purpose of calculating metrics, such as those of Table 1, is to combine them

348    later for the definition of further metrics. These are hereafter referred to as combined metrics

349    (Table 2). Several approaches have been proposed to combine geometric metrics, such as

350    metrics sum, and root mean square. The combination of metrics is done at either the local or

351    global level. For example, the index D (metric 56) combines two geometric metrics at the

352     local level ($OS_{ij}$ and $US_{ij}$) to produce a set of $D_{ij}$ values, which is then aggregated for the

353     global level. The F-measure (metric 55) combines two metrics at the global level (Precision

354     and Recall). A few more complex strategies have also been proposed for combining metrics,

355     namely clustering (CI, metric 58) and comparison of the cumulative distribution of the

356     metrics combined ($M^g$ and $M^j$, metrics 60 and 63).

357     Place table 2 near here. **<u>See Table 2 after the references</u>**.

358     Further methods are found in the literature. Most of them are essentially the same as those

359     presented in Table 1 and Table 2. They are omitted here as are ambiguously described in the

360     original publications; for example, the correspondence between objects and polygons is

361     frequently unclear. Thus, they could not be translated to the notation defined in Section 3.1.1.

362     Methods not described here are, however, potentially useful and include those found in

363     Winter (2000); Oliveira et al. (2003); Radoux and Defourny (2007); Esch et al. (2008);

364     Corcoran et al. (2010); Korting et al. (2011); Verbeeck et al. (2012); Whiteside et al. (2014);

365     Michel et al. (2015) and Mikes et al. (2015).

366     **3.1.3 Metrics use**

367     Table 1 reveals that a variety of strategies has been adopted to compare objects and polygons.

368     Specifically, often the assessment is focused on the reference data set, and thus the

369     assessment proceeds by searching the objects that may correspond to each polygon (i.e. set Y

370     is compared to set X). For example, Recall (metric 2) uses this strategy. Sometimes the

371     assessment proceeds by searching the polygons that may correspond to each object (i.e. X is

372     compared to Y). Precision (metric 1) adopts this latter strategy. The remaining strategies

373     defined in Sections 3.1.1.3 and 3.1.1.4 are less frequently adopted, namely in three specific

374     methods which calculate metrics 11-12, 40-42, and 65-66.

375     Once the strategy of comparison between objects and polygons is specified, several criteria

376     may be used to determine the correspondence between objects and polygons. For example,

377     when set Y compares to set X a simple criterion is to consider only one corresponding object

378     for each of the polygons. This object may be the one that covers the largest extent of the

379     polygon (e.g. Recall, metric 2). However, a set of different criteria can be used. For example,

380     qLoc (metric 16) views an object as corresponding to a polygon if the centroid of the polygon

381     lies inside the object or *vice versa*. As a result, several objects may be identified as

382     corresponding to a single polygon. Only the corresponding objects and polygons are used for

383     calculating the geometric metrics.

384     Most of the metrics presented in Table 1 and Table 2 are based on proportions of overlapping

385     area. For example, Precision (metric 1) is based on the calculation of the proportion of the

386     area that each object has in common with the corresponding polygon. On the other hand,

387     some metrics are based on the distance between centroids. For example, qLoc (metric 16) is

388     based on the distance between the centroid of each of the polygons to that of the

389     corresponding objects. Metrics that focus on area are often referred to as area coincidence-

390     based or area-based metrics. The metrics that focus on position are often referred to as

391     boundary coincidence-based, location-based, or position-based metrics (Cheng et al., 2014;

392     Clinton et al., 2010; Montaghi et al., 2013; Whiteside et al., 2014; Winter, 2000).

393     A substantial proportion of the metrics detect either under-segmentation or over-segmentation

394     error. This may be unexpected as commonly a balanced result is desired, but it informs on

395     what type of error dominates. This may be used, for example, to parameterize a segmentation

396     algorithm. For this reason, normally metrics that detect and measure under- or over-

397     segmentation error are calculated separately, but combined later (Table 2) to provide a

398     complementary view on image segmentation accuracy. Moreover, area-based metrics and

17

399  position-based metrics are sometimes combined to provide a comprehensive assessment of

400  image segmentation accuracy from a geometric point of view (Möller et al., 2013). The

401  combined metrics are typically the outcome of an image segmentation accuracy assessment

402  based on a geometric approach. The possible values of these metrics are in the range between

403  0 and 1, and they may be used to rank a set of image segmentation outputs based on their

404  expected suitability for image classification. To assist in the comparison of all metrics

405  presented here, the metrics of Table 1 and Table 2 are grouped in Table 3 by type of error

406  measured (over- and/or under-segmentation) and geometric feature considered (area and/or

407  position).

408  Table 3. Geometric metrics of tables 1 and 2 grouped by type of error measured (over-

409  segmentation and/or under-segmentation) and type of metric (area-based and/or position-

410  based). Combined metrics of table 2 are in bold.

| Type of metric | Type of error | | | | | |
|---|---|---|---|---|---|---|
| | Over-segmentation | | Under-segmentation | | Over- and under-segmentation | |
| Area-based | Recall | (2) | Precision | (1) | M | (5) |
| | uM | (3) | oM | (4) | AFI | (10) |
| | $LRE(x_i,y_j)_p$ | (12) | $LRE(y_j,x_i)_p$ | (11) | $d_{sym}$ | (13) |
| | RAsub | (17) | E | (14) | SimSize | (15) |
| | countOver | (26) | RAsuper | (18) | $G_s$ | (21) |
| | BsO | (30) | PI | (22) | $F_{ij}$ | (23) |
| | OS | (34) | countUnder | (27) | $m_2$ | (24) |
| | ED | (35) | $A_j$ | (29) | qr | (25) |
| | FG | (36) | LP | (31) | SH | (37) |
| | NSR | (39) | EP | (32) | SOA | (50) |
| | $O^R$ | (40) | US | (33) | MOA | (51) |
| | OE | (45) | PSE | (38) | OI2 | (54) |
| | OS2 | (48) | $O^F$ | (41) | **F** | (55) |
| | OSE | (52) | CE | (44) | **D** | (56) |
| | | | US2 | (47) | **BCE** | (57) |
| | | | TSI | (49) | **ED2** | (59) |
| | | | USE | (53) | **ADI** | (61) |
| | | | | | **ED3** | (62) |
| | | | | | **SEI** | (64) |
| | | | | | **BCA($x_i,y_j)_p$** | (65) |

| | | | | | | |
|---|---|---|---|---|---|---|
| | | | | | **BCA** | (66) |
| Position-based | User's BPA | (6) | Prod.'s BPA | (7) | qLoc | (16) |
| | C' | (8) | O' | (8) | RPsub | (19) |
| | $P^R$ | (42) | $P^F$ | (43) | RPsuper | (20) |
| | | | | | modD(b) | (28) |
| | | | | | PDI | (46) |
| Area- and position-based | | | | | **CI** | (58) |
| | | | | | $M^g$ | (60) |
| | | | | | $M^j$ | (63) |

411

## 3.2 Non-geometric methods

A small number of non-geometric methods have been proposed (Table 4). Typically, this

category of methods does not require an overlay operation between a polygonal reference

data set and the image segmentation output under evaluation as they need not to be spatially

coincident. Polygons may not even be used. The requirement common to all non-geometric

methods is that the land cover class(es) associated with the objects are known. Note that non-

geometric methods are not able to explicitly inform on which type of error, under- or over-

segmentation, predominates.

Table 4. Non-geometric methods for supervised assessment of image segmentation accuracy.
All metrics detect under- and over-segmentation error.

| Reference | Focus of the method | Polygons needed[a] |
|---|---|---|
| Wang et al. (2004) | Objects' content (spectral separability of classes using the Bhattacharyya distance). | No |
| Laliberte and Rango (2009) | Classifier (Decision trees classification accuracy and Gini index). | No |
| Anders et al. (2011) | Objects' content (difference among objects and polygons on the frequency distribution of characterizing topographic attributes). | Yes |
| Yang et al. (2017) | Classifier (classification uncertainty) | No |

[a] The reference data set used is required in the form of polygons

423    Non-geometric methods essentially follow two approaches to assess the accuracy of image

424    segmentation. The first approach focuses on the content of the objects. Anders et al. (2011)

425    compared the content of objects and polygons using the frequency distribution of their

426    topographic attributes such as slope angle while mapping geomorphological features. Smaller

427    differences between frequency distributions calculated from objects and polygons of the same

428    geomorphological feature type indicated greater segmentation accuracy. However, most of

429    the non-geometric methods dispense with polygons and only require objects with known

430    spectral and thematic content. These objects may be represented in the spectral space used in

431    the segmentation analysis where the objects of different land cover classes are desirable to lie

432    in different regions so that later a classifier can allocate them to the correct class. The

433    separability of the objects in the spectral space as a function of the land cover classes they

434    represent is regarded as indicative of segmentation accuracy, and this can be assessed based

435    on, for example, the Bhattacharyya distance (Fukunaga, 1990). This is possibly the most used

436    non-geometric method (Li et al., 2015; Radoux and Defourny, 2008; Wang et al., 2004; Xun

437    and Wang, 2015).

438    The second approach used in non-geometric methods assesses image segmentation using a

439    classifier. Specifically, a series of preliminary classifications are undertaken with a set of

440    image segmentation outputs, and the classifier is used to rank the segmentations based on

441    their suitability for image classification. For example, a sample of the objects of the image

442    segmentation under evaluation can be used to train a decision tree, and the impurity of the

443    terminal nodes can be regarded as indicative of classification success; large accuracy of

444    image segmentation is expected to be related to low node impurity (Laliberte and Rango,

445    2009). Most often, however, traditional estimators of classification accuracy such as overall

446    accuracy are used (Laliberte and Rango, 2009; Smith, 2010). Thus, the classifier suggests

447    which of a set of segmentation outputs affords the largest classification accuracy. In this case,

448    samples of the objects constructed can be used for training and testing a classifier by means

449    of out-of-bag estimate or cross-validation (Laliberte and Rango, 2009; Smith, 2010).

450    Classification uncertainty rather than accuracy can also be used. If a fuzzy classifier is

451    employed, the way in which the probability of class membership is partitioned between the

452    classes can be used to calculate classification uncertainty, for example based on entropy

453    measures. Segmentation accuracy may be viewed as negatively related to the magnitude of

454    classification uncertainty (Yang et al., 2017).

455    The second approach of non-geometric image segmentation accuracy assessment, especially

456    when classification accuracy expressed by traditional estimators such as overall accuracy is

457    considered, may appear similar to traditional classification accuracy assessment, but they are

458    different things. The former uses the training sample to assess the accuracy of the preliminary

459    classifications while the latter assesses the quality of the final mapping product and requires

460    an independent testing sample. Sometimes traditional classification accuracy assessment is

461    nevertheless used to assess indirectly image segmentation accuracy (e.g. Kim et al., 2009; Li

462    et al., 2011). When used, the focus is typically on a comparison among the accuracy values of

463    a set of final classifications (Foody, 2009, 2004), with each produced with different image

464    segmentation outputs. The differences are caused not only by the image segmentations used,

465    but the entire approach to image classification. This may be well suited for applications

466    focused on the final mapping products, but implies possibly impractical labour and resources

467    such as multiple testing samples.

# 4 Selecting a method

The selection of a method to assess the accuracy of image segmentation is a complex decision, and here it is suggested to tackle that decision from two central perspectives: the application in-hand, and the pros and cons of the methods. These issues should be considered holistically although discussed separately hereafter.

## 4.1 Application in-hand

The purpose of the application in-hand should be considered, and there are two main situations. First, the applications are focused on just a fraction of the classes in which a landscape may be categorized. These applications use image segmentation primarily for object recognition and extraction, such as buildings and trees in urban environments (e.g. Belgiu and Drăguţ, 2014; Sebari and He, 2013). The desired characteristics of the objects are likely to be geometric, such as position and shape. Several methods may be appropriate, such as shape error (metric 37); the segmentation output indicated as optimal will in principle be formed by objects that most resemble the desired shapes represented in the reference data set. Alternatively, the relative overlapping areas between objects and polygons may be maximised. This strategy may benefit from area-based metrics designed for object recognition, such as SEI (metric 64).

The second main situation corresponds to wall-to-wall land cover classification and mapping (e.g Bisquert et al. 2015; Strasser and Lang 2015). In this case, the geometric properties of the objects may be considered important as in the first situation described above, and hence geometric methods may be used. However, the thematic information associated with the objects is commonly regarded as more important than the geometrical representation. In this context, an output that enables the maximisation of the area under analysis correctly

22

491 represented in the final map is preferred. Geometric methods can still be used, and area-based

492 methods may be appropriate, which will in principle suggest as optimal the segmentation

493 output formed by objects that represent the largest amount of area of the corresponding

494 polygons. This gives the classification stage the opportunity of maximising the area correctly

495 classified and thus the overall accuracy of the map. Non-geometric methods can also be used

496 (Table 4). There is less experience in the use of this category of methods, but it is potentially

497 useful when the geometry of the objects does not have to meet predefined requirements.

498 An intermediate situation is also possible in that both the geometric and thematic properties

499 of the objects are regarded as important. In this case, methods that combine different

500 approaches for the accuracy assessment may be used, for example focused on the relative

501 position and area of overlap between objects and polygons (Möller et al., 2013, 2007).

502 However, there is no need to select just one method, and assembling multiple methods is a

503 valid option (Clinton et al., 2010). Different methods, including geometric and non-geometric

504 methods, can be used together to address all the specific properties of the objects considered

505 as relevant as long as the set of methods used fits the purpose of the application in-hand.

506 Another relevant aspect of the application in-hand is the relative importance of under- and

507 over-segmentation error. Image segmentation is typically conducted to trade-off and

508 minimize under- and over-segmentation error, but over-segmentation may be needed to

509 address conveniently the problem under analysis. Specifically, small objects, sometimes

510 called primitive objects (Dronova, 2015), may be needed for modelling complex classes that

511 are not directly related to spectral data, such as habitats (Strasser and Lang, 2015). The final

512 land cover classes can be delineated later, for example, based on knowledge-driven semantic

513 rules (Gu et al., 2017). If no primitive objects are needed, and the border of the final land

514 over classes to be mapped are pursued in a segmentation analysis, it may be desirable

515 nevertheless to recognize that under- and over-segmentation error are not always equally

516 serious, especially if the application is interested more on the thematic rather than the

517 geometric properties of the objects. Multiple authors have expressed their preference for

518 over- rather than under-segmentation error as the latter is associated with relatively small

519 classification accuracy (Gao et al., 2011; Hirata and Takahashi, 2011; Lobo, 1997; Wang et

520 al., 2004). Under-segmentation error produces objects that correspond to more than one class

521 on the ground and thus may represent an important origin of misclassification or land cover

522 map error. Therefore, using methods able to inform on the level of over- and under-

523 segmentation error may be convenient, such as that proposed by Möller et al. (2013).

524 The third and last aspect highlighted here relates to the potential importance of thematic

525 errors associated with under-segmentation error. That is, the impact of under-segmentation

526 error may depend on the classes associated with under-segmented objects. This is because the

527 needs of the individual users may vary greatly in their sensitivity to misclassifications as a

528 function of the classes involved (Bontemps et al., 2012; Comber et al., 2012). Traditionally,

529 supervised methods consider all under-segmentation errors as equally serious, but under-

530 segmentation errors can in fact be weighted as a function of the classes involved. This is the

531 situation with the geometric method proposed by Costa et al. (2015) (metric 63) and non-

532 geometric methods that use a classifier to perform a preliminary series of classifications,

533 whose results can be expressed through weighted estimators of classification accuracy, such

534 as the Stehman's (1999) map value V.

## 4.2 Methods' pros and cons

536 A consideration of the potential implications associated with the approach of the assessment

537 is advisable. Non-geometric methods do not require geo-registered reference data, which may

538   be very practical, but are unable to explicitly inform on which type of segmentation error

539   predominates. That information may be useful for guiding the definition of segmentation

540   settings. If this limitation is undesirable, a geometric method suited to detecting segmentation

541   error explicitly should be preferred. However, the need of defining criteria of correspondence

542   between objects and polygons should be considered carefully as it impacts on the accuracy

543   assessment. The geometric methods proposed by Yang et al. (2015) (SEI, metric 64), Su and

544   Zhang (2017) (OSE, metric 52), and Möller et al. (2013) ($M^g$, metric 60) pay particular

545   attention to this issue.

546   Quantitative comparisons of different methods should be undertaken. Several comparative

547   studies dedicated to geometric methods have been published (Clinton et al., 2010; Räsänen et

548   al., 2013; Whiteside et al., 2014; Yang et al., 2015), and some of them (e.g. Clinton et al.,

549   2010; Verbeeck et al., 2012) observed that different methods can indicate very different

550   segmentation outputs as optimal. Thus, special attention should be given to potential bias of

551   the methods. For example, Radoux and Defourny (2008) found that spectral separability

552   measures used in non-geometric methods may be insensitive to under-segmentation error, and

553   thus indicate a segmentation as optimal while notably under-segmented; Witharana and

554   Civco (2014) found that the sensitivity of Euclidean distance 2 (ED2, metric 59) to the

555   accuracy of the objects depends of the scale of the analysis.

556   Finally, it should be noted that estimated bias in image segmentation accuracy assessment is

557   not caused merely by unsuitable choice of methods or their potential flaws, but the protocol

558   used for their implementation. Typically, some reference data are available for a sample of

559   the entire area to be mapped, and thus limited data are used to infer an accuracy estimate to

560   represent the entire area. Therefore, the nature of sampling is an issue that will impact on the

561   results of an image segmentation accuracy assessment. The reference data must be acquired

25

562  using a probability sampling design, which must incorporate a randomization component that

563  has a non-zero probability of selection for each object into the sample. Consideration of

564  general sampling and statistical principles for defining samples is recommended (Olofsson et

565  al., 2014; Stehman and Czaplewski, 1998).

## 5   Discussion

### 5.1   Current status

568  Image segmentation accuracy assessment appears to be in a relatively early stage of

569  maturation in land cover mapping applications. Often no information on the assessment

570  produced is given, and qualitative assessment based on visual interpretation is widely used.

571  This situation may be a result of several factors. For example, the lack of a solid background

572  in image segmentation accuracy assessment and reliable recommendations for method

573  selection may be a motivation for neglecting a quantitative accuracy assessment. Another

574  factor may be related to the difficulty of implementing most of the methods proposed in the

575  literature. Many analysts of remote sensing data depend on standard software and have no

576  resources or expertise to implement new methods. This may also be a reason why comparison

577  among methods has been addressed in a relatively small number of studies. There are some

578  initiatives to implement supervised methods and make them available to the public (Mikes et

579  al. 2015), but further work should be done in this respect. Clinton et al. (2010), Montaghi et

580  al. (2013), Eisank et al. (2014), and Novelli et al. (2017) provide additional information on

581  how to access software that includes supervised methods for image segmentation accuracy

582  assessment.

583  Supervised methods were reviewed here and grouped into two categories: geometric and non-

584  geometric methods. The former includes numerous area-based methods (Table 3), and many

26

585    of them are similar. This is the case of area($x_i \cap y_j$)/area($y_j$), which appears in metrics 1, 18,

586    33, and 47. Winter (2000) demonstrated that only seven metrics are possible to derive from

587    an area-based approach if they are free of dimension, normalized, and symmetric (i.e. there is

588    a single and mutual correspondence between objects and polygons). However, several

589    correspondence criteria and strategies of comparison between objects and polygons can be

590    specified, and thus the number of area-based metrics can proliferate. This is essentially the

591    case of metrics 1, 18, 33 and 47, which are calculated with different criteria of

592    correspondence between objects and polygons ( $X'_j$, $\tilde{Y}_i$, $Y'_i$, and $Yc_i \cup Yd_i$, respectively).

593    The ways the local metric values are used to produce a global accuracy value also vary.

594    These apparently slight differences may, however, impact substantially on the assessment as

595    different calculations are involved.


596    Selecting an appropriate method for image segmentation accuracy assessment is not obvious.

597    The pros and cons of the potential methods, such as ease of use and bias, should be taken into

598    account. However, it is noted that there is often neither a right nor wrong method. The

599    suitability of a method will ultimately depend on how it fits with the application in-hand.


600    **5.2   Research needs**

601    Quantitative studies similar to Clinton et al. (2010) and Witharana and Civco (2014) should

602    be done to exhaustively test and compare the supervised methods used in the remote sensing

603    community. Non-geometric methods should be inspected as they have been neglected in

604    quantitative studies. Moreover, the studies should be conducted under different contexts that

605    may represent different types of applications, such as object recognition, and wall-to-wall

606    mapping. Critically, research to address the relationship between segmentation and

607 classification accuracies is required, as often relations were not simple (Belgiu and Drăguţ,

608 2014; Costa et al., 2017; Räsänen et al., 2013; Verbeeck et al., 2012).

609 Finally, the concept of over- and under-segmentation error should be revisited. Commonly, as

610 in this paper, segmentation error is defined relative to the reference data used, and thus the

611 concept lacks theoretical robustness. For example using reference data representing final land

612 cover classes to be mapped or primitive objects impacts on the results. Primitive objects have

613 a more spectral rather than thematic significance, and this may influence the assessment,

614 including the selection of the assessment approach, supervised or unsupervised. However,

615 theory and concepts related to object-based image analysis are generally incipient (Blaschke

616 et al., 2014; Ma et al., 2017), and comparing supervised and unsupervised methods which

617 often focus on thematic and primitive objects, respectively, has not received much attention.

## 6   Conclusions

619 Accuracy assessment is an important component of an image segmentation analysis, but is

620 not mature. It has been much undertaken through visual inspection possibly for practical

621 reasons while many quantitative approaches and methods have been proposed. Most often

622 these methods are supervised and focus on the geometry of the objects constructed and

623 polygons taken as reference data. However, other approaches may be used. The spectrum of

624 methods available is large, and it is difficult to select consciously suitable methods for

625 particular applications. There are at least three important questions that should be asked

626 during the selection of supervised methods for image segmentation accuracy assessment: (i)

627 the goal of the application; (ii) the relative importance of under- and over-segmentation error

628 (including a possible varying sensitivity to thematic issues associated to under-segmentation);

629 and (iii) the pros and cons of the methods. Answering these questions will help select suitable

630 methods, but further research is needed to improve the standards of image segmentation

631 accuracy assessment, otherwise there is the risk of using methods unsuitable or sub-optimal

632 for the application in-hand.

## References

639 Abeyta, A., Franklin, J., 1998. The accuracy of vegetation stand boundaries derived from
640     image segmentation in a desert environment. Photogramm. Eng. Remote Sensing 64,
641     59–66.

642 Anders, N.S., Seijmonsbergen, A.C., Bouten, W., 2011. Segmentation optimization and
643     stratified object-based analysis for semi-automated geomorphological mapping. Remote
644     Sens. Environ. 115, 2976–2985. doi:10.1016/j.rse.2011.05.007

645 Basaeed, E., Bhaskar, H., Hill, P., Al-Mualla, M., Bull, D., 2016. A supervised hierarchical
646     segmentation of remote-sensing images using a committee of multi-scale convolutional
647     neural networks. Int. J. Remote Sens. 37, 1671–1691.
648     doi:10.1080/01431161.2016.1159745

649 Beauchemin, M., Thomson, K.P.B., Edwards, G., 1998. On the Hausdorff distance used for
650     the evaluation of segmentation results. Can. J. Remote Sens. 24, 3–8.
651     doi:10.1080/07038992.1998.10874685

652 Belgiu, M., Drăguţ, L., 2014. Comparing supervised and unsupervised multiresolution
653     segmentation approaches for extracting buildings from very high resolution imagery.
654     ISPRS J. Photogramm. Remote Sens. 96, 67–75. doi:10.1016/j.isprsjprs.2014.07.002

655 Bisquert, M., Bégué, A., Deshayes, M., 2015. Object-based delineation of homogeneous
656     landscape units at regional scale based on MODIS time series. Int. J. Appl. Earth Obs.
657     Geoinf. 37, 72–82. doi:10.1016/j.jag.2014.10.004

658 Blaschke, T., 2010. Object based image analysis for remote sensing. ISPRS J. Photogramm.
659     Remote Sens. 65, 2–16. doi:10.1016/j.isprsjprs.2009.06.004

660 Blaschke, T., Hay, G.J., Kelly, M., Lang, S., Hofmann, P., Addink, E.A., Queiroz Feitosa, R.,
661     van der Meer, F., van der Werff, H., van Coillie, F., Tiede, D., 2014. Geographic Object-
662     Based Image Analysis – Towards a new paradigm. ISPRS J. Photogramm. Remote Sens.
663     87, 180–191. doi:10.1016/j.isprsjprs.2013.09.014

664 Bontemps, S., Herold, M., Kooistra, L., van Groenestijn, A., Hartley, A., Arino, O., Moreau,
665     I., Defourny, P., 2012. Revisiting land cover observation to address the needs of the
666     climate modeling community. Biogeosciences 9, 2145–2157. doi:10.5194/bg-9-2145-
667     2012

668 Bradley, B.A., 2014. Remote detection of invasive plants: A review of spectral, textural and
669     phenological approaches. Biol. Invasions 16, 1411–1425. doi:10.1007/s10530-013-
670     0578-9

671 Cardoso, J.S., Corte-Real, L., 2005. Toward a generic evaluation of image segmentation.
672     IEEE Trans. Image Process. 14, 1773–1782. doi:10.1109/TIP.2005.854491

673 Carleer, A.P., Debeir, O., Wolff, E., 2005. Assessment of very high spatial resolution satellite
674     image segmentations. Photogramm. Eng. Remote Sensing 71, 1285–1294.
675     doi:10.14358/PERS.71.11.1285

676 Cheng, J., Bo, Y., Zhu, Y., Ji, X., 2014. A novel method for assessing the segmentation
677     quality of high-spatial resolution remote-sensing images. Int. J. Remote Sens. 35, 3816–
678     3839. doi:10.1080/01431161.2014.919678

679 Clinton, N., Holt, A., Scarborough, J., Yan, L., Gong, P., 2010. Accuracy assessment
680     measures for object-based image segmentation goodness. Photogramm. Eng. Remote
681     Sensing 76, 289–299.

682 Coillie, F.M.B. Van, Gardin, S., Anseel, F., 2014. Variability of operator performance in
683     remote-sensing image interpretation: the importance of human and external factors. Int.
684     J. Remote Sens. 35, 754–778. doi:10.1080/01431161.2013.873152

685 Comber, A., Fisher, P., Brunsdon, C., Khmag, A., 2012. Spatial analysis of remote sensing
686     image classification accuracy. Remote Sens. Environ. 127, 237–246.
687     doi:10.1016/j.rse.2012.09.005

688 Corcoran, P., Winstanley, A., Mooney, P., 2010. Segmentation performance evaluation for
689     object-based remotely sensed image analysis. Int. J. Remote Sens. 31, 617–645.
690     doi:10.1080/01431160902894475

691 Costa, G.A.O.P., Feitosa, R.Q., Cazes, T.B., Feijó, B., 2008. Genetic adaptation of
692     segmentation parameters, in: Blaschke, T., Lang, S., Hay, G.J. (Eds.), Object-Based
693     Image Analysis. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 679–695.
694     doi:10.1007/978-3-540-77058-9_37

695 Costa, H., Foody, G.M., Boyd, D.S., 2015. Integrating user needs on misclassification error
696     sensitivity into image segmentation quality assessment. Photogramm. Eng. Remote
697     Sensing 81, 451–459. doi:10.14358/PERS.81.6.451

698    Costa, H., Foody, G.M., Boyd, D.S., 2017. Using mixed objects in the training of object-
699        based image classifications. Remote Sens. Environ. 190, 188–197.
700        doi:10.1016/j.rse.2016.12.017

701    Crevier, D., 2008. Image segmentation algorithm development using ground truth image data
702        sets. Comput. Vis. Image Underst. 112, 143–159. doi:10.1016/j.cviu.2008.02.002

703    Drăguţ, L., Csillik, O., Eisank, C., Tiede, D., 2014. Automated parameterisation for multi-
704        scale image segmentation on multiple layers. ISPRS J. Photogramm. Remote Sens. 88,
705        119–127. doi:10.1016/j.isprsjprs.2013.11.018

706    Drăguţ, L., Tiede, D., Levick, S.R., 2010. ESP: a tool to estimate scale parameter for
707        multiresolution image segmentation of remotely sensed data. Int. J. Geogr. Inf. Sci. 24,
708        859–871. doi:10.1080/13658810903174803

709    Dronova, I., 2015. Object-based image analysis in wetland research: A review. Remote Sens.
710        7, 6380–6413. doi:10.3390/rs70506380

711    Eisank, C., Smith, M., Hillier, J., 2014. Assessment of multiresolution segmentation for
712        delimiting drumlins in digital elevation models. Geomorphology 214, 452–464.
713        doi:10.1016/j.geomorph.2014.02.028

714    Esch, T., Thiel, M., Bock, M., Roth, A., Dech, S., 2008. Improvement of image segmentation
715        accuracy based on multiscale optimization procedure. IEEE Geosci. Remote Sens. Lett.
716        5, 463–467. doi:10.1109/LGRS.2008.919622

717    Feitosa, R.Q., Ferreira, R.S., Almeida, C.M., Camargo, F.F., Costa, G.A.O.P., 2010.
718        Similarity metrics for genetic adaptation of segmentation parameters, in: 3rd
719        International Conference on Geographic Object-Based Image Analysis (GEOBIA 2010).
720        The International Archives of the Photogrammetry, Remote Sensing and Spatial
721        Information Sciences, Ghent.

722    Feizizadeh, B., Blaschke, T., Tiede, D., Moghaddam, M.H.R., 2017. Evaluating fuzzy
723        operators of an object-based image analysis for detecting landslides and their changes.
724        Geomorphology 293, Part A, 240–254.
725        doi:https://doi.org/10.1016/j.geomorph.2017.06.002

726    Foody, G.M., 2009. Classification accuracy comparison: Hypothesis tests and the use of
727        confidence intervals in evaluations of difference, equivalence and non-inferiority.
728        Remote Sens. Environ. 113, 1658–1663. doi:10.1016/j.rse.2009.03.014

729    Foody, G.M., 2004. Thematic map comparison: Evaluating the statistical significance of
730        differences in classification accuracy. Photogramm. Eng. Remote Sensing 70, 627–633.
731        doi:10.14358/PERS.70.5.627

732    Fukunaga, M., 1990. Introduction to statistical pattern recognition, 2nd ed. Academic Press,
733        San Diego.

734    Gao, Y., Mas, J.F., Kerle, N., Pacheco, J.A.N., 2011. Optimal region growing segmentation
735        and its effect on classification accuracy. Int. J. Remote Sens. 32, 3747–3763.

736        doi:10.1080/01431161003777189

737    Gu, H., Li, H., Yan, L., Liu, Z., Blaschke, T., Soergel, U., 2017. An Object-Based Semantic
738        Classification Method for High Resolution Remote Sensing Imagery Using Ontology.
739        Remote Sens. . doi:10.3390/rs9040329

740    Hirata, Y., Takahashi, T., 2011. Image segmentation and classification of Landsat Thematic
741        Mapper data using a sampling approach for forest cover assessment. Can. J. For. Res.
742        41, 35–43. doi:10.1139/X10-130

743    Hultquist, C., Chen, G., Zhao, K., 2014. A comparison of Gaussian process regression,
744        random forests and support vector regression for burn severity assessment in diseased
745        forests. Remote Sens. Lett. 5, 723–732. doi:10.1080/2150704X.2014.963733

746    Janssen, L.L.F., Molenaar, M., 1995. Terrain objects, their dynamics and their monitoring by
747        the integration of GIS and remote sensing. IEEE Trans. Geosci. Remote Sens. 33, 749–
748        758. doi:10.1109/36.387590

749    Kim, M., Madden, M., Warner, T.A., 2009. Forest type mapping using object-specific texture
750        measures from multispectral Ikonos Imagery: Segmentation quality and image
751        classification issues. Photogramm. Eng. Remote Sensing 75, 819–829.

752    Korting, T.S., Dutra, L.V., Fonseca, L.M.G., 2011. A resegmentation approach for detecting
753        rectangular objects in high-resolution imagery. IEEE Geosci. Remote Sens. Lett. 8, 621–
754        625. doi:10.1109/LGRS.2010.2098389

755    Laliberte, A.S., Rango, A., 2009. Texture and scale in object-based analysis of subdecimeter
756        resolution unmanned aerial vehicle (UAV) imagery. IEEE Trans. Geosci. Remote Sens.
757        47, 1–10. doi:10.1109/TGRS.2008.2009355

758    Lang, S., Albrecht, F., Kienberger, S., Tiede, D., 2010. Object validity for operational tasks
759        in a policy context. J. Spat. Sci. 55, 9–22. doi:10.1080/14498596.2010.487639

760    Lang, S., Kienberger, S., Tiede, D., Hagenlocher, M., Pernkopf, L., 2014. Geons – domain-
761        specific regionalization of space. Cartogr. Geogr. Inf. Sci. 41, 214–226.
762        doi:10.1080/15230406.2014.902755

763    Levine, M.D., Nazif, A.M., 1982. An experimental rule based system for testing low level
764        segmentation strategies, in: Preston, K., Uhr, L. (Eds.), Multicomputers and Image
765        Processing: Algorithms and Programs. Academic Press, New York, pp. 149–160.

766    Li, D., Ke, Y., Gong, H., Li, X., 2015. Object-based urban tree species classification using bi-
767        temporal WorldView-2 and WorldView-3 images. Remote Sens. 7, 16917–16937.
768        doi:10.3390/rs71215861

769    Li, P., Guo, J., Song, B., Xiao, X., 2011. A multilevel hierarchical image segmentation
770        method for urban impervious surface mapping using very high resolution imagery. IEEE
771        J. Sel. Top. Appl. Earth Obs. Remote Sens. 4, 103–116.
772        doi:10.1109/JSTARS.2010.2074186

773    Liu, D., Xia, F., 2010. Assessing object-based classification: Advantages and limitations.
774       Remote Sens. Lett. 1, 187–194. doi:10.1080/01431161003743173

775    Liu, J., Li, P., Wang, X., 2015. A new segmentation method for very high resolution imagery
776       using spectral and morphological information. ISPRS J. Photogramm. Remote Sens.
777       101, 145–162. doi:10.1016/j.isprsjprs.2014.11.009

778    Liu, Y., Bian, L., Meng, Y., Wang, H., Zhang, S., Yang, Y., Shao, X., Wang, B., 2012.
779       Discrepancy measures for selecting optimal combination of parameter values in object-
780       based image analysis. ISPRS J. Photogramm. Remote Sens. 68, 144–156.
781       doi:10.1016/j.isprsjprs.2012.01.007

782    Lobo, A., 1997. Image segmentation and discriminant analysis for the identification of land
783       cover units in ecology. IEEE Trans. Geosci. Remote Sens. 35, 1136–1145.
784       doi:10.1109/36.628781

785    Lucieer, A., Stein, A., 2002. Existential uncertainty of spatial objects segmented from
786       satellite sensor imagery. Geosci. Remote Sensing, IEEE Trans. 40, 2518–2521.
787       doi:10.1109/TGRS.2002.805072

788    Ma, L., Li, M., Ma, X., Cheng, L., Du, P., Liu, Y., 2017. A review of supervised object-based
789       land-cover image classification. ISPRS J. Photogramm. Remote Sens. 130, 277–293.
790       doi:10.1016/j.isprsjprs.2017.06.001

791    Marpu, P.R., Neubert, M., Herold, H., Niemeyer, I., 2010. Enhanced evaluation of image
792       segmentation results. J. Spat. Sci. 55, 55–68. doi:10.1080/14498596.2010.487850

793    Martin, D.R., 2003. An empirical approach to grouping and segmentation. ECS Department,
794       University of California.

795    Matikainen, L., Karila, K., Hyyppä, J., Litkey, P., Puttonen, E., & Ahokas, E. (2017). Object-
796       based analysis of multispectral airborne laser scanner data for land cover classification
797       and map updating. *ISPRS Journal of Photogrammetry and Remote Sensing*, *128*, 298–
798       313. http://doi.org/https://doi.org/10.1016/j.isprsjprs.2017.04.005

799    Michel, J., Youssefi, D., Grizonnet, M., 2015. Stable mean-shift algorithm and its application
800       to the segmentation of arbitrarily large remote sensing images. IEEE Trans. Geosci.
801       Remote Sens. 53, 952–964. doi:10.1109/TGRS.2014.2330857

802    Mikes, S., Haindl, M., Scarpa, G., Gaetano, R., 2015. Benchmarking of remote sensing
803       segmentation methods. IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens. 8, 2240–2248.
804       doi:10.1109/JSTARS.2015.2416656

805    Möller, M., Birger, J., Gidudu, A., Gläßer, C., 2013. A framework for the geometric accuracy
806       assessment of classified objects. Int. J. Remote Sens. 34, 8685–8698.
807       doi:10.1080/01431161.2013.845319

808    Möller, M., Lymburner, L., Volk, M., 2007. The comparison index: A tool for assessing the
809       accuracy of image segmentation. Int. J. Appl. Earth Obs. Geoinf. 9, 311–321.
810       doi:10.1016/j.jag.2006.10.002

811    Montaghi, A., Larsen, R., Greve, M.H., 2013. Accuracy assessment measures for image
812        segmentation goodness of the Land Parcel Identification System (LPIS) in Denmark.
813        Remote Sens. Lett. 4, 946–955. doi:10.1080/2150704X.2013.817709

814    Neubert, M., Herold, H., Meinel, G., 2008. Assessing image segmentation quality –
815        Concepts, methods and application, in: Blaschke, T., Lang, S., Hay, G. (Eds.), Object-
816        Based Image Analysis. Springer Berlin Heidelberg, pp. 769–784.

817    Novelli, A., Aguilar, M., Aguilar, F., Nemmaoui, A., Tarantino, E., 2017. AssesSeg—A
818        Command Line Tool to Quantify Image Segmentation Quality: A Test Carried Out in
819        Southern Spain from Satellite Imagery. Remote Sens. 9, 40. doi:10.3390/rs9010040

820    Oliveira, J., Formaggio, A., Epiphanio, J., Luiz, A., 2003. Index for the Evaluation of
821        Segmentation (IAVAS): An application to agriculture. Mapp. Sci. Remote Sens. 40,
822        155–169. doi:10.2747/0749-3878.40.3.155

823    Olofsson, P., Foody, G.M., Herold, M., Stehman, S. V., Woodcock, C.E., Wulder, M.A.,
824        2014. Good practices for estimating area and assessing accuracy of land change. Remote
825        Sens. Environ. 148, 42–57. doi:10.1016/j.rse.2014.02.015

826    Persello, C., Bruzzone, L., 2010. A novel protocol for accuracy assessment in classification
827        of very high resolution images. IEEE Trans. Geosci. Remote Sens. 48, 1232–1244.
828        doi:10.1109/TGRS.2009.2029570

829    Qi, Z., Yeh, A.G.-O., Li, X., Zhang, X., 2015. A three-component method for timely
830        detection of land cover changes using polarimetric SAR images. ISPRS J. Photogramm.
831        Remote Sens. 107, 3–21. doi:10.1016/j.isprsjprs.2015.02.004

832    Radoux, J., Defourny, P., 2008. Quality assessment of segmentation results devoted to object-
833        based classification, in: Blaschke, T., Lang, S., Hay, G.J. (Eds.), Object-Based Image
834        Analysis. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 257–271.
835        doi:10.1007/978-3-540-77058-9_14

836    Radoux, J., Defourny, P., 2007. A quantitative assessment of boundaries in automated forest
837        stand delineation using very high resolution imagery. Remote Sens. Environ. 110, 468–
838        475. doi:10.1016/j.rse.2007.02.031

839    Räsänen, A., Rusanen, A., Kuitunen, M., Lensu, A., 2013. What makes segmentation good?
840        A case study in boreal forest habitat mapping. Int. J. Remote Sens. 34, 8603–8627.
841        doi:10.1080/01431161.2013.845318

842    Robson, B.A., Nuth, C., Dahl, S.O., Hölbling, D., Strozzi, T., Nielsen, P.R., 2015. Automated
843        classification of debris-covered glaciers combining optical, SAR and topographic data in
844        an object-based environment. Remote Sens. Environ. 170, 372–387.
845        doi:10.1016/j.rse.2015.10.001

846    Sebari, I., He, D.-C., 2013. Automatic fuzzy object-based analysis of VHSR images for urban
847        objects extraction. ISPRS J. Photogramm. Remote Sens. 79, 171–184.
848        doi:10.1016/j.isprsjprs.2013.02.006

849 Smith, A., 2010. Image segmentation scale parameter optimization and land cover
850      classification using the Random Forest algorithm. J. Spat. Sci. 55, 69–79.
851      doi:10.1080/14498596.2010.487851

852 Stehman, S. V., 1999. Comparing thematic maps based on map value. Int. J. Remote Sens.
853      20, 2347–2366. doi:10.1080/014311699212065

854 Stehman, S. V., Czaplewski, R.L., 1998. Design and analysis for thematic map accuracy
855      assessment. Remote Sens. Environ. 64, 331–344. doi:10.1016/S0034-4257(98)00010-8

856 Strasser, T., Lang, S., 2015. Object-based class modelling for multi-scale riparian forest
857      habitat mapping. Int. J. Appl. Earth Obs. Geoinf. 37, 29–37.
858      doi:10.1016/j.jag.2014.10.002

859 Su, T., Zhang, S., 2017. Local and global evaluation for remote sensing image segmentation.
860      ISPRS J. Photogramm. Remote Sens. 130, 256–276. doi:10.1016/j.isprsjprs.2017.06.003

861 Tian, J., Chen, D.-M., 2007. Optimization in multi-scale segmentation of high-resolution
862      satellite images for artificial feature recognition. Int. J. Remote Sens. 28, 4625–4644.
863      doi:10.1080/01431160701241746

864 Van Coillie, F.M.B., Verbeke, L.P.C., De Wulf, R.R., 2008. Semi-automated forest stand
865      delineation using wavelet based segmentation of very high resolution optical imagery,
866      in: Object-Based Image Analysis: Spatial Concepts for Knowledge-Driven Remote
867      Sensing Applications. pp. 237–256. doi:10.1007/978-3-540-77058-9_13

868 Van Rijsbergen, C.J., 1979. Information retrieval. Butterworth-Heinemann, London.

869 Verbeeck, K., Hermy, M., Van Orshoven, J., 2012. External geo-information in the
870      segmentation of VHR imagery improves the detection of imperviousness in urban
871      neighborhoods. Int. J. Appl. Earth Obs. Geoinf. 18, 428–435.
872      doi:10.1016/j.jag.2012.03.015

873 Wang, L., Sousa, W.P., Gong, P., 2004. Integration of object-based and pixel-based
874      classification for mapping mangroves with IKONOS imagery. Int. J. Remote Sens. 25,
875      5655–5668. doi:10.1080/014311602331291215

876 Weidner, U., 2008. Contribution to the assessment of segmentation quality for remote sensing
877      applications. Int. Arch. Photogramm. Remote Sens. Spat. Inf. Sci. 37, 479–484.

878 Whiteside, T.G., Maier, S.W., Boggs, G.S., 2014. Area-based and location-based validation
879      of classified image objects. Int. J. Appl. Earth Obs. Geoinf. 28, 117–130.
880      doi:10.1016/j.jag.2013.11.009

881 Winter, S., 2000. Location similarity of regions. ISPRS J. Photogramm. Remote Sens. 55,
882      189–200. doi:10.1016/S0924-2716(00)00019-8

883 Witharana, C., Civco, D.L., 2014. Optimizing multi-resolution segmentation scale using
884      empirical methods: Exploring the sensitivity of the supervised discrepancy measure
885      Euclidean distance 2 (ED2). ISPRS J. Photogramm. Remote Sens. 87, 108–121.

886        doi:10.1016/j.isprsjprs.2013.11.006

887  Witharana, C., Civco, D.L., Meyer, T.H., 2014. Evaluation of data fusion and image
888        segmentation in earth observation based rapid mapping workflows. ISPRS J.
889        Photogramm. Remote Sens. 87, 1–18. doi:10.1016/j.isprsjprs.2013.10.005

890  Xun, L., Wang, L., 2015. An object-based SVM method incorporating optimal segmentation
891        scale estimation using Bhattacharyya Distance for mapping salt cedar (Tamarisk spp.)
892        with QuickBird imagery. GIScience Remote Sens. 52, 257–273.
893        doi:10.1080/15481603.2015.1026049

894  Yang, J., He, Y., Caspersen, J.P., Jones, T., 2017. Delineating Individual Tree Crowns in an
895        Uneven-Aged, Mixed Broadleaf Forest Using Multispectral Watershed Segmentation
896        and Multiscale Fitting. IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens. 10, 1390–1401.
897        doi:10.1109/JSTARS.2016.2638822

898  Yang, J., He, Y., Caspersen, J., Jones, T., 2015. A discrepancy measure for segmentation
899        evaluation from the perspective of object recognition. ISPRS J. Photogramm. Remote
900        Sens. 101, 186–192. doi:10.1016/j.isprsjprs.2014.12.015

901  Yang, J., Li, P., He, Y., 2014. A multi-band approach to unsupervised scale parameter
902        selection for multi-scale image segmentation. ISPRS J. Photogramm. Remote Sens. 94,
903        13–24. doi:10.1016/j.isprsjprs.2014.04.008

904  Yi, L., Zhang, G., Wu, Z., 2012. A scale-synthesis method for high spatial resolution remote
905        sensing image segmentation. IEEE Trans. Geosci. Remote Sens. 50, 4062–4070.
906        doi:10.1109/TGRS.2012.2187789

907  Zhan, Q., Molenaar, M., Tempfli, K., Shi, W., 2005. Quality assessment for geo-spatial
908        objects derived from remotely sensed data. Int. J. Remote Sens. 26, 2953–2974.
909        doi:10.1080/01431160500057764

910  Zhang, H., Fritts, J.E., Goldman, S.A., 2008. Image segmentation evaluation: A survey of
911        unsupervised methods. Comput. Vis. Image Underst. 110, 260–280.
912        doi:10.1016/j.cviu.2007.08.003

913  Zhang, X., Feng, X., Xiao, P., He, G., Zhu, L., 2015a. Segmentation quality evaluation using
914        region-based precision and recall measures for remote sensing images. ISPRS J.
915        Photogramm. Remote Sens. 102, 73–84. doi:10.1016/j.isprsjprs.2015.01.009

916  Zhang, X., Xiao, P., Feng, X., Feng, L., Ye, N., 2015b. Toward evaluating multiscale
917        segmentations of high spatial resolution remote sensing images. IEEE Trans. Geosci.
918        Remote Sens. 53, 3694–3706. doi:10.1109/TGRS.2014.2381632

919  Zhang, X., Xiao, P., Feng, X., Wang, J., Wang, Z., 2014. Hybrid region merging method for
920        segmentation of high-resolution remote sensing images. ISPRS J. Photogramm. Remote
921        Sens. 98, 19–28. doi:10.1016/j.isprsjprs.2014.09.011

922  Zhang, Y.J., 1996. A survey on evaluation methods for image segmentation. Pattern
923        Recognit. 29, 1335–1346. doi:10.1016/0031-3203(95)00169-7

924

925

926    LIST OF FIGURE CAPTIONS

927    Figure 1. Geometric comparison between an object and polygon based on the overlapping

928    area (shaded area) and/or distance between centroids (dashed arrow).

929    Figure 2. Sets X, Y, and S: (a) reference set X, (b) segmentation Y, and (c) intersection

930    $S = X \cap Y$. In (c) yellow denotes one-to-one, blue denotes one-to-many, and pink denotes

931    many-to-many (Section 3.1.1.3).

932    Figure 3. Comparison between X and Y of Figure 2: (a) four potential objects (dashed lines)

933    corresponding to polygon $x_1$ (grey background) when Y compares to X; (b) one potential

934    polygon (dashed line) corresponding to object $y_1$ (grey background) when X compares to Y.

Table 1. Geometric metrics for supervised assessment of image segmentation accuracy. All metrics are numbered and ordered chronologically. The type of metric and segmentation error are identified in columns Typ. and Err. while the minimum, maximum, and optimal values of the metrics are identified in columns Min., Max., and Opt. The subscripts of the metrics' name indicate local accuracy assessment (see notes on the corresponding global metric), and global metrics have no subscripts.

| Metric | Reference | Typ.[a] | Err.[b] | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|---|
| (1) $\text{Precision}_{ij} = \dfrac{\text{area}(x_i \cap y_j)}{\text{area}(y_j)}, x_i \in X'_j$ | Van Rijsbergen (1979) and Zhang et al. (2015a). | AB | U | 0 | 1 | 1 | Global metric Precision is the weighted mean of all Precision$_{ij}$ using area(y$_j$) as weights. |
| (2) $\text{Recall}_{ij} = \dfrac{\text{area}(x_i \cap y_j)}{\text{area}(x_i)}, y_j \in Y'_i$ | Van Rijsbergen (1979) and Zhang et al. (2015a). | AB | O | 0 | 1 | 1 | Global metric Recall is the weighted mean of all Recall$_{ij}$ using area(x$_i$) as weights. |
| (3) $\text{underMerging}_{ij} = \dfrac{\text{area}(x_i) - \text{area}(x_i \cap y_j)}{\text{area}(x_i)}, y_j$ | Levine and Nazif (1982) and Clinton et al. (2010). | AB | O | 0 | 0.5 | 0 | Global metric underMerging can be the mean of all underMerging$_{ij}$. |
| (4) $\text{overMerging}_{ij} = \dfrac{\text{area}(y_j) - \text{area}(x_i \cap y_j)}{\text{area}(x_i)}, y_j \in$ | Levine and Nazif (1982) and Clinton et al. (2010). | AB | U | 0 | 0.5 | 0 | Global metric overMerging can be the mean of all overMerging$_{ij}$. |

| Metric | Reference | Typ.[a] | Err.[b] | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|---|
| (5) $M_{ij} = \sqrt{\dfrac{area(x_i \cap y_j)^2}{area(x_i) \times area(y_j)}}, y_j \in Y_i'$ | Janssen and Molenaar (1995) and Feitosa et al. (2010). | AB | UO | 0 | 1 | 1 | Match (M). Global metric M is the mean of all $M_{ij}$ values. |
| (6) User's BPA= proportion of boundary length defined in segmentation with corresponding real boundaries | Abeyta and Franklin (1998) | PB | O | 0 | 1 | 1 | Boundary positional accuracy (BPA). Boundary length are estimated based on point-type data collected via line intersect sampling. Boundaries defined in segmentation that fell within ε (epsilon) tolerances (spatial error bounds) of surveyed boundaries are considered accurate. |
| (7) Producers's BPA= proportion of real boundary length with corresponding boundaries defined in segmentation | Abeyta and Franklin (1998) | PB | U | 0 | 1 | 1 | Boundary positional accuracy (BPA). Boundary length are estimated based on point-type data collected via line intersect sampling. Boundaries defined in segmentation that fell within ε (epsilon) tolerances (spatial error bounds) of surveyed boundaries are considered accurate. |

| Metric | Reference | Typ.[a] | Err.[b] | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|---|
| (8) $C' = 1 - \dfrac{size\left(\bigcup_i vertex(y_j)\right) - size\left(\bigcup_i (dist(vertex(y_j), vertex(x_i)))\right)}{size\left(\bigcup_i vertex(x_i)\right)}, x_i \in X$ | Beauchemin et al. (1998). | PB | O | 0 | | 0 | dist() represents the partial directed Hausdorff distance, which calculates the fraction of vertexes of the objects of Y that are each within a distance of some vertex of the polygons of X. |
| (9) $O' = 1 - \dfrac{size\left(\bigcup_i (dist(vertex(x_i), vertex(y_j)))\right)}{size\left(\bigcup_i vertex(x_i)\right)}, y_j \in Y$ | Beauchemin et al. (1998). | PB | U | 0 | 1 | 0 | dist() represents the partial directed Hausdorff distance, which calculates the fraction of vertexes of the polygons of X that are each within a distance of some vertex of the objects of Y. |
| (10) $AFI_{ij} = \dfrac{area(x_i) - area(y_j)}{area(x_i)}, y_j \in Y'_i$ | Lucieer and Stein (2002) and Clinton et al. (2010). | AB | UO | | | 0 | Area fit index (AFI). Global metric AFI is the mean of all $AFI_{ij}$ values. AFI<0 and AFI>0 indicate under- and over-segmentation. |
| (11) $LRE(y_j, x_i)_p = \dfrac{size(\neg x_i \cap y_j)}{size(y_j)}, x_i \in Xa_p \wedge y_j \in Ya_p$ | Martin (2003) and Zhang et al. (2015a). | AB | U | 0 | 1 | 0 | Local refinement error (LRE). This metric was not proposed to be aggregated for the entire segmentation output (see metric 57 in Table 2). |
| (12) $LRE(x_i, y_j)_p = \dfrac{size(x_i \cap \neg y_j)}{size(x_i)}, x_i \in Xa_p \wedge y_j \in Ya_p$ | Martin (2003) and Zhang et al. (2015a). | AB | O | 0 | 1 | 0 | Local refinement error (LRE). This metric was not proposed to be aggregated for the entire segmentation output (see metric 57 in Table 2). |

| Metric | Reference | Typ.[a] | Err.[b] | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|---|
| (13) $d_{sym}$ = minimal number of pixels that must be removed from both X and Y so that they are identical in the remaining pixels. | Cardoso and Corte-Real (2005). | AB | UO | 0 | 1 | 0 | $d_{sym}$ is normalized to 0-1 by dividing by $l$-1. $d'_{sym} = 1 - d_{sym}$ in Zhang et al. (2015a). |
| (14) $E_{ij} = \dfrac{area(y_j) - area(x_i \cap y_j)}{area(y_j)} \times 100, x_i \in X'_j$ | Carleer et al. (2005). | AB | U | 0 | 50 | 0 | Global metric E is the weighted mean of all $E_{ij}$, using $area(y_j)$ as weights. A refinement of E is also presented in Carleer et al. (2005). |
| (15) $SimSize_{ij} = \dfrac{min(area(x_i), area(y_j))}{max(area(x_i), area(y_j))}, y_j \in Y_i^*$ | Zhan et al. (2005) and Clinton et al. (2010). | AB | UO | 0 | 1 | 1 | Global metric SimSize can be the mean of all $SimSize_{ij}$. |
| (16) $qLoc_{ij} = dist(centroid(x_i), centroid(y_j)), y_j \in Y_i^*$ | Zhan et al. (2005) and Clinton et al. (2010). | PB | UO | 0 | | 0 | dist() represents Euclidean distance. Global metric qLoc can be the mean of all $qLoc_{ij}$. |
| (17) $RAsub_{ij} = \dfrac{area(x_i \cap y_j)}{area(x_i)}, y_j \in \tilde{Y}_i$ | Möller et al. (2007) and Clinton et al. (2010). | AB | O | 0 | 1 | 1 | Relative area (RA). This metric was not proposed to be aggregated for the whole segmentation output (see metric 58 in Table 2). |

| Metric | Reference | Typ.[a] | Err.[b] | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|---|
| (18) $\mathrm{RA\,super_{ij}} = \dfrac{\mathrm{area}(x_i \cap y_j)}{\mathrm{area}(y_j)}, y_j \in \tilde{Y}_i$ | Möller et al. (2007) and Clinton et al. (2010). | AB | U | 0 | 1 | 1 | Relative area(RA). This metric was not proposed to be aggregated for the whole segmentation output (see metric 58 in Table 2). |
| (19) $\mathrm{RPsub}_{ij} = \mathrm{dist\,(centroid\,}(x_i), \mathrm{centroid\,}(y_j)), y_j \in \tilde{Y}_i$ | Möller et al. (2007) and Clinton et al. (2010). | PB | UO | 0 | | 0 | Relative position (RP). This metric was not proposed to be aggregated for the whole segmentation output (see metric 58 in Table 2). |
| (20) $\mathrm{RP\,super_{ij}} = \dfrac{\mathrm{dist(centroid}(x_i), \mathrm{centroid}(y_j))}{\max(\mathrm{RPsub}_{ij})}, y_j \in Y_i^*$ | Möller et al. (2007) and Clinton et al. (2010). | PB | UO | 0 | 1 | 0 | Relative position (RP). dist() represents Euclidean distance. This metric was not proposed to be aggregated for the whole segmentation output (see metric 58 in Table 2). |
| (21) $G_s = \dfrac{\sum_i \sum_j \mathrm{area}(x_i \cap y_j)}{\mathrm{area}(X) \times e^{\sum_i \sum_j \left(\mathrm{area}(x_i \cup y_j) - \mathrm{area}(x_i \cap y_j)\right) / \mathrm{area}(X)}}, y_j \in Yc_i$ | Tian and Chen (2007). | AB | UO | 0 | 1 | 1 | |
| (22) $\mathrm{PI}_{ij} = \dfrac{\mathrm{area}(x_i \cap y_j)^2}{\mathrm{area}(y_j) \times \mathrm{area}(x_i)}, y_j \in \tilde{Y}_i$ | Coillie et al. (2008). | AB | U | 0 | 1 | 1 | Purity Index (PI). Global metric PI is the mean of all summed $\mathrm{PI}_{ij}$ over all $x_i$. |

43

| Metric | Reference | Typ.[a] | Err.[b] | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|---|
| (23) $F_{ij} = \dfrac{\text{area}(y_j) + \text{area}(x_i) - 2 \times \text{area}(y_j \cap x_i)}{\text{area}(y_j)}, x_i \in X'_j$ | Costa et al. (2008). | AB | UO | 0 | | 0 | Fitness function (F). Global metric F is the mean of all summed $F_{ij}$ over all $y_j$. |
| (24) $m_{2_{ij}} = \dfrac{\text{area}(y_j \cap x_i)}{\text{area}(y_j \cup x_i)}, x_i \in X''_j$ | Crevier (2008) and Yi et al. (2012). | AB | UO | 0 | 1 | 1 | Global metric $m_2$ is the sum of all $m_{2ij}$. |
| (25) $qr_{ij} = 1 - \dfrac{\text{area}(x_i \cap y_j)}{\text{area}(x_i \cup y_j)}, y_j \in Y_i^*$ | Weidner (2008) and Clinton et al. (2010). | AB | UO | 0 | 1 | 0 | Quality rate (qr). Global metric qr can be the mean of all $qr_{ij}$ |
| (26) $\text{countOver} = \text{size}(X), \dfrac{\text{area}(y_j)}{\text{area}(x_i)} < 1 \wedge AFI_{ij} > 0 \wedge y_j \in Y_i^*$ | Clinton et al. (2010). | AB | O | 0 | size(x) | 0 | $AFI_{ij} = \dfrac{\text{area}(x_i) - \text{area}(y_j)}{\text{area}(x_i)}, y_j \in Y'_i$ (see metric 10). |
| (27) $\text{countUnder} = \text{size}(X), \dfrac{\text{area}(y_j)}{\text{area}(x_i)} = 1 \wedge AFI_{ij} < 0 \wedge y_j \in Y_i^*$ | Clinton et al. (2010). | AB | U | 0 | size(x) | 0 | $AFI_{ij} = \dfrac{\text{area}(x_i) - \text{area}(y_j)}{\text{area}(x_i)}, y_j \in Y'_i$ (see metric 10). |

| Metric | Reference | Typ.[a] | Err.[b] | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|---|
| (28) $\text{modD(b)}_i$ = mean Euclidean distance between each vertex of $x_i$ and the closest vertex in every $y_j \in Y_i^*$ | Clinton et al. (2010). | PB | UO | 0 | | 0 | Global metric modD(b) can be the mean of all $\text{modD(b)}_i$. |
| (29) $A_j = \dfrac{\max(\text{area}(c_i \cap y_j))}{\text{area}(y_j)}, c_i \in \tilde{C}_j$ | Liu and Xia (2010). | AB | U | 0 | 1 | 1 | Segmentation accuracy (A). Global metric A is the weighted mean of all $A_j$ using area($y_j$) as weights. |
| (30) $\text{BsO}_i = \max\left(\dfrac{\text{area}(y_j) - \text{area}(\neg x_i \cap y_j)}{\text{area}(x_i)}\right) \times 100, y_j \in Yf_i$ | Marpu et al. (2010). | AB | O | 0 | 100 | 100 | Biggest sub-object (BsO). Global BsO can be descriptive statistics of all $\text{BsO}_i$ (e.g. quartiles). |
| (31) $\text{LP}_i = \dfrac{\text{area}(x_i) - \sum_j \text{area}(x_i \cap y_j)}{\text{area}(x_i)} \times 100, y_j \in Yf_i$ | Marpu et al. (2010). | AB | U | 0 | 100 | 0 | Lost pixels (LP). Global LP can be descriptive statistics of all $\text{LP}_i$ (e.g. quartiles). |
| (32) $\text{EP}_{ij} = \dfrac{\text{area}(y_j) - \text{area}(x_i \cap y_j)}{\text{area}(x_i)} \times 100, y_j \in Yf_i$ | Marpu et al. (2010). | AB | U | 0 | 100 | 0 | Extra pixels (EP). Global EP can be descriptive statistics of all summed $\text{EP}_{ij}$ over all $x_i$ (e.g. quartiles). |

| Metric | Reference | Typ.[a] | Err.[b] | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|---|
| (33) $US_{ij} = 1 - \dfrac{area(x_i \cap y_j)}{area(y_j)}, y_j \in Y_i'$ | Persello and Bruzzone (2010) and Clinton et al. (2010). | AB | U | 0 | 1 | 0 | undersegmentation error (US). Global metric US can be the mean of all $US_{ij}$. Clinton et al. (2010) consider subset $Y_i^*$. |
| (34) $OS_{ij} = 1 - \dfrac{area(x_i \cap y_j)}{area(x_i)}, y_j \in Y_i'$ | Persello and Bruzzone (2010) and Clinton et al. (2010). | AB | O | 0 | 1 | 0 | oversegmentation error (OS). Global metric OS can be the mean of all $OS_{ij}$. Clinton et al. (2010) consider subset $Y_i^*$. |
| (35) $ED_{ij} = 1 - \dfrac{perim(x_i) \cap perim(y_j)}{perim(x_i)}, y_j \in Y_i'$ | Persello and Bruzzone (2010). | AB | O | 0 | 1 | 0 | Edge location (ED). Global metric ED can be the mean of all $ED_{ij}$. |
| (36) $FG_i = \dfrac{size(\tilde{Y}_i) - 1}{area(x_i) - 1}$ | Persello and Bruzzone (2010). | AB | O | 0 | 1 | 0 | Fragmentation error (FG). Global metric FG can be the mean of all $FG_i$. |
| (37) $SH_{ij} = \left| sf(x_i) - sf(y_j) \right|, y_j \in Y_i'$ | Persello and Bruzzone (2010). | AB | UO | 0 | | 0 | Shape error (SH). $|\cdot|$ denotes the absolute value of '·' and $sf(\cdot)$ denotes a shape factor of '·' such as compactness and sphericity. Global metric SH can be the mean of all $SH_{ij}$. |

| Metric | Reference | Typ.[a] | Err.[b] | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|---|
| (38) $PSE_{ij} = \dfrac{area(\neg x_i \cap y_j)}{area(x_i)}, y_j \in Yc_i \cup Yd_i$ | Liu et al. (2012). | AB | U | 0 | | 0 | Potential segmentation error (PSE). Global metric PSE is the weighted mean all $PSE_{ij}$, using $area(x_j)$ as weights. A refinement of PSE is presented in Novelli et al. (2017). |
| (39) $NSR = \dfrac{\left\lvert size(X) - size\left(\bigcup_i \left(Yc_i \cup Yd_i\right)\right)\right\rvert}{size(X)}$ | Liu et al. (2012). | AB | O | 0 | | 0 | Number-of-segments ratio (NSR). $\lvert \cdot \rvert$ denotes the absolute value of '$\cdot$'. A refinement of NSR is presented in Novelli et al. (2017). |
| (40) $O^R_{ijk} = \dfrac{area(s_{ijk})}{area(x_i)}, s_{ijk} \in Sax_i \vee Sbx_i$ | Möller et al. (2013). | AB | O | 0 | 1 | 1 | This metric was not proposed to be aggregated for the whole segmentation output (see metric 60 in Table 2). |
| (41) $O^F_{ijk} = \dfrac{area(s_{ijk})}{area(y_j)}, s_{ijk} \in Say_j \vee Sby_j$ | Möller et al. (2013). | AB | U | 0 | 1 | 1 | This metric was not proposed to be aggregated for the whole segmentation output (see metric 60 in Table 2). |
| (42) $P^R_{ijk} = 1 - \dfrac{dist(centroid(s_{ijk}), centroid(x_i))}{d^x_{max}}, s_{ijk} \in Sax_i \vee Sbx_i$ | Möller et al. (2013). | PB | O | 0 | 1 | 1 | $d^x_{max} = max(dist(centroid(s_{ijk}))), s_{ijk} \in Sax_i \vee Sbx_i$. dist() represents Euclidean distance. This metric was not proposed to be aggregated for the whole segmentation output (see metric 60 in Table 2). |

| Metric | Reference | Typ.[a] | Err.[b] | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|---|
| (43) $P_{ijk}^F = 1 - \dfrac{dist(centroid(s_{ijk}), centroid(y_j))}{d_{max}^y}, s_{ijk} \in Say_j \vee Sby_j$ | Möller et al. (2013). | PB | U | 0 | 1 | 1 | $d_{max}^y = max(dist(centroid(s_{ijk}))), s_{ijk} \in Say_j \vee Sby_j$. dist() represents Euclidean distance. This metric was not proposed to be aggregated for the whole segmentation output (see metric 60 in Table 2). |
| (44) $CE_{ij} = \dfrac{area(y_i) - area(x_i \cap y_j)}{area(x_i)} \times 100, y_j \in Yb_i \cap Yc_i$ | Cheng et al. (2014). | AB | U | 0 | 50 | 0 | Commission error (CE). Global metric $CE_{overall}$ is the weighted mean of all $CE_{ij}$, using area($x_j$) as weights. |
| (45) $OE_{ij} = \dfrac{area(x_i \cap y_j)}{area(x_i)} \times 100, y_j \in \tilde{Y}_i \setminus Yb_i \cap Yc_i$ | Cheng et al. (2014). | AB | O | 0 | 50 | 0 | Omission error (OE). Global metric $OE_{overall}$ is the weighted mean of all $OE_{ij}$, using area($x_j$) as weights. |
| (46) $PDI_{ij} = dist(centroid(x_i), centroid(y_j)), y_j \in Yb_i \cup Yc_i$ | Cheng et al. (2014). | PB | UO | 0 | | 0 | Position discrepancy index (PDI). Global metric $PDI_{overall}$ is the mean of all averaged $PDI_{ij}$ over all $x_i$. |
| (47) $US2_{ij} = 1 - \dfrac{area(x_i \cap y_j)}{area(y_j)}, y_j \in Yc_i \cup Yd_i$ | Yang et al. (2014). | AB | U | 0 | 1 | 0 | Global metric US is the sum of all summed $US_{ij}$ over each $x_i$. |
| (48) $OS2_{ij} = 1 - \dfrac{area(x_i \cap y_j)}{area(x_i)}, y_j \in Yc_i \cup Yd_i$ | Yang et al. (2014). | AB | O | 0 | 1 | 0 | Global metric OS is the sum of all summed $OS_{ij}$ over each $x_i$. |

| Metric | Reference | Typ.[a] | Err.[b] | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|---|
| (49) $$\mathrm{TSI}_j = \sum_{c_i}\left(\frac{\mathrm{area}(c_i)}{\mathrm{area}(y_j)}\sum_{d_i}\left(\frac{\mathrm{area}(d_i)}{\mathrm{area}(y_j)}w_{c_i d_i}\right)\right), c_i \wedge d_i \in \tilde{C}_j$$ | Costa et al. (2015). | AB | U | 0 | 1 | 1 | Thematic similarity index (TSI). Global metric TSI is the weighted mean of all $\mathrm{TSI}_j$. using $\mathrm{area}(y_j)$ as weights. |
| (50) $$\mathrm{SOA}_{ij} = \frac{\mathrm{area}(x_i \cap y_j) \times 2}{\mathrm{area}(x_i) + \mathrm{area}(y_j)}, y_j \in \tilde{Y}_i$$ | Zhang et al. (2015b). | AB | UO | 0 | 1 | 1 | Single-scale object accuracy (SOA). This metric was not proposed to be aggregated for the whole segmentation output, but only for each $x_i$, which is $\mathrm{SOA}_i = \max(\mathrm{SOA}_{ij})$ |
| (51) $$\mathrm{MOA}_i = \max(\{\mathrm{SOA}_i\}), \mathrm{size}(\{\mathrm{SOA}_i\}) = h$$ | Zhang et al. (2015b). | AB | UO | 0 | 1 | 1 | Multiscale object accuracy (MOA). Metric developed to assess multiscale segmentation, that is, several sets Y are created ($Y_1$, $Y_2$,… $Y_h$), from which a set of $h$ metrics $\mathrm{SOA}_i$ are calculated for each $x_i$. $\mathrm{SOA}_i$ corresponds to metric 50. Global metric MOA is the weighted mean of all $\mathrm{MOA}_i$, using $\mathrm{area}(x_j)$ as weights. |
| (52) $$\mathrm{OSE}_i = \begin{cases} \dfrac{1}{1 - \dfrac{1}{\mathrm{area}(x_i)}}\left(1 - \dfrac{\mathrm{area}(x_i \cap y_j)}{\mathrm{area}(x_i)}\right)_{ij}, & y_j \in Yg_i \\ 0, & y_j \notin Yg_i \end{cases}$$ | Su and Zhang (2017). | AB | O | 0 | 1 | 0 | Over-segmentation error (OSE). Global metric OSE (called GOSE) is the weighted mean of all $\mathrm{OSE}_j$. using $\mathrm{area}(x_i)$ as weights. |

| Metric | Reference | Typ.[a] | Err.[b] | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|---|
| (53) $$USE_i = \frac{\min\left(\begin{array}{l} area(x_i), \\ \left(area(x_i) - \sum_j area(x_i \cap y_j)\right) + \\ \left(\sum_j area(x_i \cup y_j) - area(x_i)\right) \end{array}\right)}{area(x_i)}, y_j \in Yg_i$$ | Su and Zhang (2017). | AB | U | 0 | 1 | 0 | Under-segmentation error (USE). Global metric USE (called GUSE) is the weighted mean of all $USE_j$. using $area(x_i)$ as weights. |
| (54) $$OI2_i = \max\left(\frac{area(x_i \cap y_j)}{area(x_i)} \times \frac{area(x_i \cap y_j)}{area(y_j)}\right), y_j \in \tilde{Y}_i$$ | Yang et al. (2017). | AB | UO | 0 | 1 | 1 | Overlap index (OI2). This metric was not proposed to be aggregated for the whole segmentation output. |

[a] area-based (AB) or position-based (PB)

[b] under-segmentation (U), over-segmentation (O), or both (UO)

Table 2. Combined geometric metrics based on those described in Table 1. The information associated with each of the columns is presented as in Table 1. All metrics detect under-segmentation and over-segmentation error.

| Combined metric | Reference | Typ. | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|
| (55) $\text{F} - \text{measure} = \dfrac{1}{\dfrac{\alpha}{\text{Precision}} + (1-\alpha)\dfrac{1}{\text{R}}}$ | Van Rijsbergen (1979) and Zhang et al. (2015a). | AB | 0 | 1 | 1 | $\alpha=0.5$ in Zhang et al. (2015a). Further combined metrics based on Precision and recall (metrics 1-2) are found in Zhang et al. (2015a). |
| (56) $\text{D}_{ij} = \sqrt{\dfrac{\text{OS}_{ij}^2 + \text{US}_{ij}^2}{2}}$ | Levine and Nazif (1982) and Clinton et al. (2010). | AB | 0 | 1 | 0 | Index D (D). Global metric D can be the mean of all $\text{D}_{ij}$. More similar combined metrics are found in Clinton et al. (2010). See metrics 33-34. |
| (57) $\text{BCE}_p = \max\big(\text{LRE}(x_i, y_j)_p, \text{LRE}(y_j, x_i)_p\big)$ | Martin (2003) and Zhang et al. (2015a). | AB | 0 | 1 | 0 | Bidirectional consistency error (BCE). Global metric BCE is the mean of all $\text{BCE}_p$. See metrics 11-12. |

| Combined metric | Reference | Typ. | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|
| (58) $$CI = \frac{\sum_{i=1}^{k}(C_i \times A_{C_i})}{k}$$ | Möller et al. (2007). | AB and PB | 0 | 100 | 100 | Comparison index (CI). $C_i$ is the comparison class, which represents clustered and ranked object metrics of over- and under-segmentation such as RAsub and RAsuper (metrics 17-18). $C_i$ can be calculated with a clustering algorithm such as $K$-means. $A_{C_i}$ is equivalent to the proportion of $C_i$ within the reference space. |
| (59) $ED2 = \sqrt{PSE^2 + NSR^2}$ | Liu et al. (2012). | AB | 0 | | 0 | Euclidian distance 2 (ED2). See metrics 38-39. |

| Combined metric | Reference | Typ. | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|
| (60) $M^g = D^- - D^+$ | Möller et al. (2013). | AB and PB | 0 | 1 | 1 | $D^-$ and $D^+$ are the distance between the cumulative distribution functions of metrics $G_{ijk}^R$ and $G_{ijk}^F$ measured by a Kolmogorov–Smirnov test, in which the null hypothesis is that the distribution function of $G_{ijk}^R$ is not less or not greater than that of $G_{ijk}^F$, respectively. $G_{ijk}^R = \sqrt{O_{ijk}^R \times P_{ijk}^R}$ (see metrics 40 and 42) and $G_{ijk}^F = \sqrt{O_{ijk}^F \times P_{ijk}^F}$ (see metrics 41 and 43). |
| (61) $ADI_{ij} = \sqrt{OE_{ij}^2 + CE_{ij}^2}$ | Cheng et al. (2014). | AB | 0 | | 0 | Area discrepancy index (ADI). Global metric $ADI_{overall} = \sqrt{OE_{overall}^2 + CE_{overall}^2}$ (see metrics 44-45). |

| Combined metric | Reference | Typ. | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|
| (62) $$ED3_{ij} = \sqrt{\frac{(OS2_{ij})^2 + (US2_{ij})^2}{2}}$$ | Yang et al. (2014). | AB | 0 | 1 | 0 | Euclidean distance 3 (ED3). Global metric ED3 is the sum of all summed $ED3_{ij}$ over each $x_i$. See metrics 47-48. ED3 modified in Yang et al. (2015b). |
| (63) $M^j = D^- - D^+$ | Costa et al. (2015). | AB and PB | 0 | 1 | 1 | $M^j$ is analogous to $M^g$ (metric 60) and $D^-$ and $D^+$ are the distance between the cumulative distribution functions of metrics $J^F_{ijk}$ and $J^R_{ijk}$ measured by a Kolmogorov–Smirnov test, in which the null hypothesis is that the distribution function of $J^F_{ijk}$ is not less or not greater than that of $J^R_{ijk}$, respectively. $J^R_{ijk} = \sqrt{G^R_{ijk}}$ and $J^F_{ijk} = \sqrt{G^F_{ijk} \times TSI_j}$ (see metric 49 and notes of metrics 60). |
| (64) $$SEI_i = \begin{cases} ED3_{ij}, y_j \in Yc_i \cap Yd_i \\ 1 \quad , y_j \notin Yc_i \cap Yd_i \end{cases}$$ | Yang et al. (2015a). | AB | 0 | 1 | 0 | Segmentation evaluation index (SEI). Global metric SEI is the mean of all $SEI_i$. |

| Combined metric | Reference | Typ. | Min. | Max. | Opt. | Notes |
|---|---|---|---|---|---|---|
| (65) $BCA(x_i, y_j)_p = BCA(y_j, x_i)_p =$ $\min(1 - LRE(x_i, y_j)_p, 1 - LRE(y_j, x_i)_p)$ | Zhang et al. (2015b). | AB | 0 | 1 | 0 | Bidirectional consistency accuracy (BCA). This metric was not proposed to be aggregated for the global level (see metrics 11, 12, and 66). |
| (66) $BCA_p = \max(\{BCA(x_i, y_j)_p\}), size(\{BCA(x_i, y_j)_p\}) = h$ | Zhang et al. (2015b). | AB | 0 | 1 | 0 | Bidirectional consistency accuracy (BCA). Metric developed to assess multiscale segmentation, that is, several sets Y are created ($Y_1, Y_2,\ldots Y_h$), from which a set of $h$ metrics $BCA(x_i, y_j)_p$ (see metric 65) are calculated for each $z_p$. Global metric BCA is the mean of all $BCA_p$. |