

GeoPAT: A toolbox for pattern-based information retrieval from large geospatial databases

Jarosław Jasiewicz^{a,c,*}, Paweł Netzel^{b,c}, Tomasz Stepinski^c

^a*Institute of Geocology and Geoinformation, Adam Mickiewicz University in Poznan, Dziegielowa 27, 60-680 Poznan*

^b*Department of Climatology and Atmosphere Protection, University of Wrocław, Kosiby 6/8, 51-621 Wrocław, Poland*

^c*Space Informatics Lab, Department of Geography, University of Cincinnati, Cincinnati, OH 45221-0131, USA*

Abstract

Geospatial Pattern Analysis Toolbox (GeoPAT) is a collection of GRASS GIS modules for carrying out pattern-based geospatial analysis of images and other spatial datasets. The need for pattern-based analysis arises when images/rasters contain rich spatial information either because of their very high resolution or their very large spatial extent. Elementary units of pattern-based analysis are scenes – patches of surface consisting of a complex arrangement of individual pixels (patterns). GeoPAT modules implement popular GIS algorithms, such as query, overlay, and segmentation, to operate on the grid of scenes. To achieve these capabilities GeoPAT includes a library of scene signatures – compact numerical descriptors of patterns, and a library of distance functions – providing numerical means of assessing dissimilarity between scenes. Ancillary GeoPAT modules use these functions to construct a grid of scenes or to assign signatures to individual scenes having regular or irregular geometries. Thus GeoPAT combines knowledge retrieval from patterns with mapping tasks within a single integrated GIS environment. GeoPAT is designed to identify and analyze complex, highly generalized classes in spatial datasets. Examples include distinguishing between different styles of urban settlements using VHR images, delineating different landscape types in land cover maps, and mapping physiographic units from DEM. The concept of pattern-based spatial analysis is explained and the roles of all modules and functions are described. A case study example pertaining to delineation of landscape types in a subregion of NLCD is given. Performance evaluation is included to highlight GeoPAT's applicability to very large datasets. The GeoPAT toolbox is available for download from <http://sil.uc.edu/>.

Keywords: pattern analysis, query-by-example, large geospatial datasets, similarity, image classification, GRASS GIS

1. Introduction

Most spatial datasets in geosciences originate from remote sensing (RS) and are in the form of images. Therefore, there exists a significant body of literature on retrieving information from RS images (Richards, 1999). Image classification - a process of converting an image into a thematic map of semantically meaningful classes - is the most common form of spatial information retrieval from an image (Lu and Weng, 2007). An original approach to image classification utilizes a pixel-based methodology. A pixel is the smallest element of a surface, as depicted in an image, for which a

value of a color is stored. A pixel-based classification algorithm assigns class labels to individual pixels. Note that this is fundamentally different from how an analyst interprets an image by perceiving the coherence of colors on multiple scales simultaneously and assigning class labels to multi-pixel tracts on the basis of their textures or patterns. Pixel-based classification algorithms may suffer from poor performance especially if applied to very high resolution (VHR) images, where individual pixels correspond to small elements of real objects and their numerical attributes are not sufficient to recognize the class of an object, or, if applied to very large images where the goal of analysis is to retrieve generalized classes (for example, when the goal is to retrieve landscape types rather than their constituent land cover classes (Graesser et al., 2012; Niesterowicz and Stepinski, 2013; Vatsavai, 2013a; Jasiewicz et al., 2014)).

*Corresponding author

Email addresses: jarekj@amu.edu.pl (Jarosław Jasiewicz), pawel.netzel@uni.wroc.pl (Paweł Netzel), stepintz@ucmail.uc.edu (Tomasz Stepinski)

30 Object-Based Image Analysis (OBIA) was developed
31 (Blaschke, 2010; Lang, 2008) to alleviate the problems
32 associated with pixel-based classification. In OBIA
33 image is first segmented to simplify it by grouping
34 pixels into meaningful segments (called "objects")
35 which are homogeneous with respect to pixel-based
36 attributes. In the second step information is retrieved
37 by classifying objects into semantically meaningful
38 classes. OBIA algorithms get closer to the way an
39 analyst interprets an image but they still suffer
40 from a number of shortcomings (Vatsavai, 2013b).
41 First, segmentation itself is a complex and
42 computationally expensive process and there is no
43 single method that performs consistently well
44 (does not under-segment or over-segment portions
45 of an image) on different RS images. Second,
46 because objects are, by definition, homogeneous
47 segments of the surface, OBIA cannot be used to
48 classify an image into highly generalized classes.
49 For example, although OBIA can classify an image
50 into land cover classes (low-level generalization)
51 more accurately than a pixel-based classifier can,
52 it still cannot classify it into landscape types
53 (high-level generalization). In other words,
54 OBIA can utilize information about image texture
55 but not information about spatial patterns.

56 For the purpose of this paper we define a spatial
57 pattern as a perceptual structure, placement, or
58 arrangement of image objects having a geometric
59 quality. We then define texture as a structure of
60 pixels arranged quasi randomly and lacking
61 geometric quality. Thus, a single land cover
62 class in a VHR image (for example, a rooftop)
63 is characterized by texture as it appears on
64 image as a quasi random mosaic of pixels having
65 a range of colors. However, a fragment of a
66 thematic map showing an urban scene consisting
67 of a spatial arrangement of several land cover
68 classes needs to be characterized by its pattern.

69 The case for classifying an image or image-like
70 spatial dataset, for example a Digital Elevation
71 Model (DEM), on the basis of spatial patterns
72 arises in multiple disciplines where a high level
73 of generalization is desired. In RS, with VHR
74 images containing rich spatial information, the
75 use of a pattern-based classification method
76 makes it possible to distinguish between
77 different urban landscapes, for example, between
78 informal settlements, industrial/commercial
79 structures, and formal residential settlements
80 (Graesser et al., 2012; Vatsavai, 2013a). In
81 landscape ecology, it makes it possible to
82 distinguish between different landscape types
83 (Niestrowicz and Stepinski, 2013; Cardille and
84 Lambois, 2009) as well as between different
85 types of forest structures (Long et al., 2010),
86 and in geomorphology it makes it possible to
87 identify and delineate physio-

88 graphic units (Jasiewicz et al., 2014).

89 It is only recently that methodologies for
90 pattern-based information retrieval from images
91 and other raster datasets have been proposed.
92 Vatsavai (2013a) proposed a multi-instance
93 learning (MIL) scheme as a means for the
94 pattern-based classification of images. In this
95 method, an image is divided into regular grid
96 of local blocks of pixels. The data (a set of
97 all multi-dimensional attribute vectors from
98 each pixel) in each block is modeled using a
99 multivariate Gaussian distribution. The
100 distance (dissimilarity) between any two
101 blocks, and thus between the two patterns
102 contained in these blocks, is calculated as the
103 probabilistic distance between their modeled
104 Gaussian distributions using the Kullback-Leibler
105 (KL) divergence. Using supervised learning
106 based on the MIL scheme Vatsavai (2013a) and
107 Graesser et al. (2012) classified RS images
108 of several cities into formal and informal
109 neighborhoods.

110 Independently, we have proposed a general
111 approach for pattern-based information retrieval
112 from all types of geospatial datasets
113 (Jasiewicz and Stepinski, 2013a; Stepinski
114 et al., 2014). For our method to be broadly
115 applicable and computationally efficient it
116 uses an input (image, DEM etc.) that has
117 been preprocessed using a pixel-based
118 classification and thus already converted into
119 a categorical format. This categorical raster
120 is divided into a regular grid of local blocks
121 of pixels. Because the data is categorical,
122 each block can be compactly represented by a
123 histogram of categories or other attributes
124 derived from these categories. We have
125 successfully applied this methodology to
126 search for and classify land-cover patterns
127 in the National Land Cover Dataset (NLCD)
128 (Jasiewicz and Stepinski, 2013a). We have
129 also used it for an assessment of land cover
130 change over the entire United States using
131 the NLCD (Netzel and Stepinski, 2015), and
132 for the identification and delineation of
133 physiographic units using DEM data
134 (Jasiewicz et al., 2014).

135 The concept of pattern-based information
136 retrieval from geospatial datasets is at the
137 beginning of its developmental cycle. For this
138 concept to mature much more work is needed,
139 including application to many different
140 datasets in multiple contexts. In this paper
141 we present the Geospatial Pattern Analysis
142 Toolbox (GeoPAT) - a collection of GRASS
143 GIS modules that integrate the various tools
144 necessary for experimenting with pattern-
145 based information retrieval from geospatial
146 data. GeoPAT is intended as a convenient
147 platform for experimentation with the
148 pattern-based analysis of rasters including
149 rasters having giga-cell and larger sizes.
150 It integrates into the GIS system
151 procedures for pattern description, pattern
152 similarity, and the

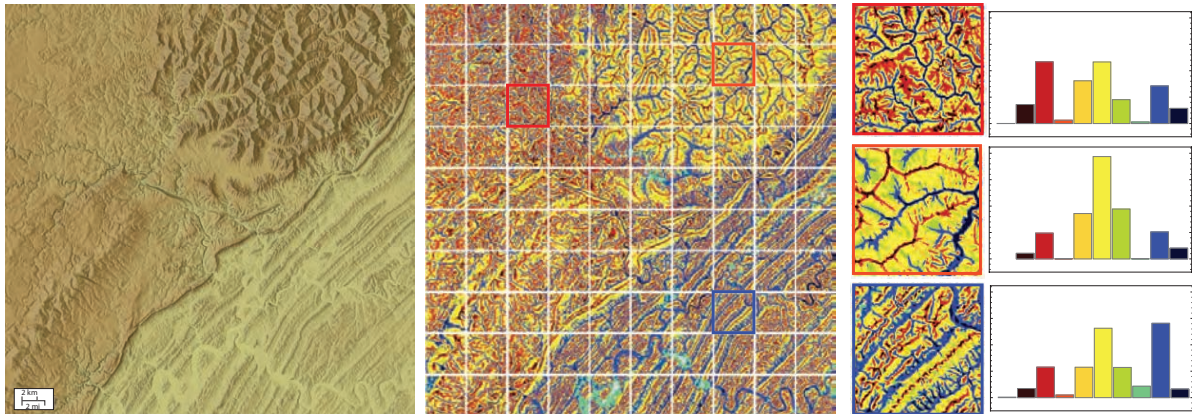


Figure 1: Example introducing a concept of pattern-based analysis of spatial datasets. (Left) Hillshade rendition of DEM over 60×60 km region. (Middle) DEM data classified into ten landform classes and divided into regular blocks. (Right) Close-ups of three sample blocks and histograms of their landform classes.

134 search and retrieval of similar patterns. These concepts
 135 were originally developed for working with natural images
 136 in the context of Content-Based Image Retrieval
 137 (CBIR) systems (Datta et al., 2008) but are now
 138 utilized by GeoPAT for the purpose of geospatial analyt-
 139 ics. Such integration allows a user to perform the stan-
 140 dard GIS tasks of mapping, map overlay, and segmen-
 141 tation on a grid of pattern-bearing blocks of pixels in a
 142 way which is already familiar (from performing simi-
 143 lar tasks on standard images). In other words, GeoPAT
 144 extends the standard GIS system by adding a new type
 145 of attribute – the pattern signature – and a new type of
 146 data query – a query-by-pattern-similarity (QBPS). This
 147 significantly lowers the cost of entry into experimenting
 148 with pattern-based information retrieval, helps to accel-
 149 erate further development of this concept, and makes
 150 possible the assessment of its utility in various domains.

151 GeoPAT modules are written in ANSI C and are
 152 designed to work within the GRASS GIS 7 (GRASS
 153 Development Team, 2012) environment. Embedding
 154 GeoPAT in GRASS has a number of advantages: (1)
 155 GRASS is an open source software available for ma-
 156 jor computing platforms, (2) GRASS is especially well-
 157 suited to work with large datasets, and (3) incorporat-
 158 ing a toolbox into an already existing, well-established
 159 environment allows for an integrated computational
 160 pipeline that provides convenience and boosts efficiency
 161 (Körting et al., 2013). GeoPAT is an actively developed
 162 solution. The core of the toolbox consists of the seven
 163 modules that compute pattern signatures and perform
 164 the GIS tasks of comparing, searching, overlaying, and
 165 segmenting the rasters on the basis of similarity between
 166 local patterns. These modules provide the basic infras-

167 tructure for pattern-based information retrieval and are
 168 not expected to be modified by a user. In addition, two
 169 libraries provide a selection of functions for extracting
 170 pattern signatures and for calculation of similarity/dis-
 171 tance between two patterns, respectively. As there are
 172 no standard means of representing spatial patterns and
 173 calculating a measure of similarity between them, we
 174 expect users to add to those libraries as they experiment
 175 with different datasets.

176 The rest of this paper is organized as follows: Sec-
 177 tion 2 presents an overview of our toolbox architec-
 178 ture. Section 3 describes the most important functions
 179 in the shared libraries and section 4 describes the seven
 180 core geoprocessing modules. A case study (section 5)
 181 presents an example on how GeoPAT modules can be
 182 utilized to perform regionalization of land cover pat-
 183 terns into landscape types using either unsupervised or
 184 supervised approaches. Section 6 gives an assessment
 185 of the computational performance of the GeoPAT mod-
 186 ules and section 7 contains our discussion and conclu-
 187 sions.

188 2. Software architecture

189 As an introduction to GeoPAT we first give an illus-
 190 tration of the basic idea behind the pattern-based anal-
 191 ysis of geospatial data. For this we use a DEM with
 192 30 m resolution. The left panel in Fig. 1 shows a hill-
 193 shade rendition of a 2000×2000 cell DEM (we reserve
 194 the term pixel for images and use the more general term
 195 *cell* for all raster datasets). The entire spatial extent
 196 of the data is referred to as a *region*. Three clearly dis-
 197 tinct physiographic units are observed in this region and

198 one of the goals of pattern-based analysis is to delineate 250
199 these units. In the preprocessing step, which is not a 251
200 part of GeoPAT, DEM cells are classified into ten land- 252
201 form classes using the geomorphons method (Jasiewicz 253
202 and Stepinski, 2013b). The result of this classification is 254
203 shown in the central panel of Fig. 1 with different colors 255
204 indicating different landforms. This panel also shows a 256
205 division of the region into a regular grid of blocks, each 257
206 block containing a large number of cells forming a local, 258
207 block-bounded pattern of landforms. A block is a 259
208 particular example of a *scene*; in general, we refer to 260
209 any subregion of the entire region as a scene. 261

210 A grid of regular scenes (as in the middle panel of 262
211 Fig. 1) is referred to as a *grid-of-scenes*. GeoPAT per- 263
212 forms GIS operations on the grid-of-scenes in the same 264
213 way as the standard GIS system performs similar opera- 265
214 tions on the grid of cells. Thus, for example, to delin- 266
215 eate the three physiographic units as seen in the sam- 267
216 ple DEM GeoPAT will classify the scenes in a way that 268
217 is analogous to how a standard pixel-based algorithm 269
218 would delineate different landform classes. Significant 270
219 technical differences in performing these operations on 271
220 scenes vs. cells stem from differences in mathematical 272
221 representations of patterns vs. numbers, and from dif- 273
222 ferences in the definitions of a distance between patterns 274
223 vs. distance between vectors. 275

224 Close-ups of three sample scenes, labeled by red, or- 276
225 ange, and blue frames are shown in the right panel of 277
226 Fig. 1 with their corresponding histograms of landform 278
227 classes. GeoPAT uses histograms as concise representa- 279
228 tions of patterns. Note that the three scenes, each rep- 280
229 resenting a different physiographic unit and exhibiting 281
230 a different pattern of landform classes, happen to have 282
231 different histograms of classes. However, in general, vi- 283
232 sually different patterns may have similar histograms of 284
233 classes. This is why GeoPAT uses more advanced his- 285
234 tograms that encapsulate not only the composition of a 286
235 pattern (the relative abundance of classes) but also its 287
236 configuration (spatial arrangement of clumps – contigu- 288
237 ous groups of same-class cells). 289

238 In general, the input to GeoPAT is a categorical raster 290
239 (classified original spatial dataset, for example, an im- 291
240 age, DEM, etc.) or a set of co-registered categorical 292
241 rasters. Additional rasters are needed for some tasks, 293
242 such as, for example, a change detection task, or to pro- 294
243 vide ancillary information for description of local pat- 295
244 terns (see section 3.1). The raster’s region is subdivided 296
245 into a regular grid (grid-of-scenes) having cells (referred 297
246 to as *s-cells*) with the size equal to or larger than the size 298
247 of the raster cells. Each *s-cell* is the center of a square 299
248 scene containing a local pattern made up of cells with 300
249 different class labels. It also stores a concise description

of this pattern which is referred to as a *signature*. The
size of the scene must be equal to or larger than the size
of the *s-cell*. This allows GeoPAT to work with over-
lapping scenes. If the size of the scene is equal to the
size of the *s-cell* the scenes don’t overlap (as in the case
shown in Fig. 1). If the size of the scene is larger than
the size of the *s-cell* the scenes overlap. In the example
given in Fig. 1 the size of the cell is 30 m and the size of
the *s-cell* is the same as the size of the scene and equal
to 6000 m.

The core of the GeoPAT toolbox consists of three
modules designed for extracting scene signatures from
categorical data, as well as four additional modules for
performing geoprocessing tasks on the grid-of-scenes.
GeoPAT implements three different signature extraction
modules: *p.sig.points*, *p.sig.polygons* and *p.sig.grid*.
This is because some geoprocessing tasks require a de-
scription of scenes not restricted to those defined by the
grid-of-scenes. For example, a search task requires a
comparison of scenes in a grid-of-scenes with a scene
(a query) defined over a region not aligned with a grid,
and a segment classification task requires calculating
signatures from irregularly-shaped scenes. The role of
geoprocessing modules (*p.sim.distmatrix*, *p.sim.search*,
p.sim.compare, *p.sim.segment*) is to perform geopro-
cessing tasks on scenes generated and described by the
signature extraction modules. The names given to the
modules adhere to the following convention: *p.* stands
for pattern, *sig.* stands for signature, and *sim.* stands for
similarity.

In addition to modules, GeoPAT provides two li-
braries of functions. The first library implements dif-
ferent methods of extracting a signature from a scene.
Functions in this library work with signature extraction
modules. The second library implements different dis-
tance measures between signatures. Functions in this
library work with geoprocessing modules. We expect
that users may want to add their own functions to both
libraries. The overall software architecture of GeoPAT
is shown in Fig. 2.

3. Library functions

Library functions implement concepts which facili-
tate working with spatial patterns in a quantitative fash-
ion. In the domain of geoscience working quantitatively
with spatial patterns has been addressed in the fields of
remote sensing (Datcu et al., 2003; Daschiel and Datcu,
2005; Li and Narayanan, 2004; Shyu et al., 2007), land-
scape ecology (Cain and Riitters, 1997; Long et al.,
2010; Cardille and Lambois, 2009; Dilts et al., 2010),
and cartography (Pontius, 2002; Remmel and Csillag,

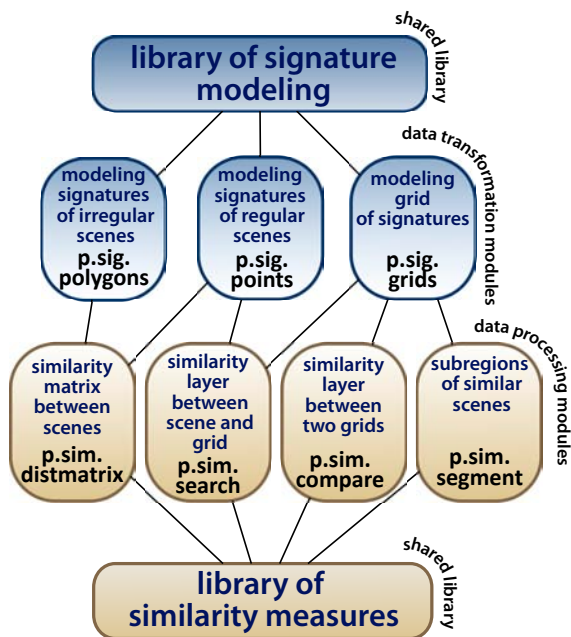


Figure 2: Architecture of GeoPAT toolbox

2006). As GeoPAT uses categorical rasters, only previous approaches developed in landscape ecology and cartography are potentially directly relevant. Because we are interested in pattern similarity measures that are rotationally invariant (for example, a scene and the same scene rotated by 90 degrees must be measured as identical), cartographic approaches, which focus on comparing different maps of the same region for consistency and accuracy, are not directly relevant.

In landscape ecology categorical patterns are described using landscape metrics (Haines-Young and Chopping, 1996; McGarigal et al., 2002; Uuemaa et al., 2009; Steiniger and Hay, 2009) which are rotationally invariant measures of compositional and configurational aspects of a scene. A collection of various landscape metrics forms an attribute vector which potentially can be used as scene signature. Several studies (Long et al., 2010; Kupfer et al., 2012; Cardille and Lambois, 2009; Cardille et al., 2012) used landscape metrics-based attribute vectors and the Euclidean distance to calculate similarities between mostly binary (forest/no forest) scenes, but the validity of such an approach has not been demonstrated. Our own experience with using landscape metrics for assessing the similarity between scenes is negative. We have identified a number of issues for using landscape metrics in GeoPAT including the selection of metrics (this can be overcome by

data reduction using PCA (Cushman et al., 2008)), the proper way to normalize metrics, and properly weighting the contribution of composition vs. configuration to the overall similarity value.

Following the principles established in the field of Content-Based Image Retrieval (CBIR) (Gevers and Smeulders, 2004; Datta et al., 2008; Lew et al., 2006) – a non-geoscience domain where the issue of similarity between two rasters (natural images) has been studied extensively – GeoPAT calculates a signature as a (possibly multi-dimensional) histogram of a pattern "primitive features." Primitive features are simple local elements of a pattern. For example, the cell's class is a primitive feature. A combination of classes of two neighboring cells is another example of a primitive feature. Many other such features could be designed. There is no generally preferred choice of primitive features; patterns in different datasets may be best encapsulated by different features. GeoPAT implements several popular methods of representing pattern by a histogram of primitive features, but it is expected that users may want to add their own.

In GeoPAT a similarity between two scenes is calculated as a similarity between two histograms, each representing a pattern contained in its respective scene. Choosing the most appropriate similarity function is largely an empirical decision which depends on the dataset and on the choice of primitive features. GeoPAT implements several histogram similarity functions, but, as with primitive features, we expect users to add their own. Cha (2007) provides a comprehensive review of histogram similarity functions.

3.1. Signature functions

Signature functions define primitive features that characterize a local pattern bounded by an extent of a scene. From among many possible signatures we describe three which are already implemented in GeoPAT.

crossproduct – This method calculates signature as a k -dimensional histogram using k primitive features assigned to each cell. Examples of such features include cell class, the size of the clump to which the cell belongs, the shape of the clump and its spatial orientation (Williams and Wentz, 2008). Because all features must be categorical (so the histogram can be formed), numerical features need to be categorized. For example, clump sizes need to be categorized into size categories from the smallest to the largest. The number of bins in the crossproduct histogram is $N_1 \times N_2 \times \dots \times N_k$, where N_i is the number of categories of i -th feature. Fig. 3A shows schematically a construction of crossproduct histogram from two features, cell class ($N_1=4$ categories depicted

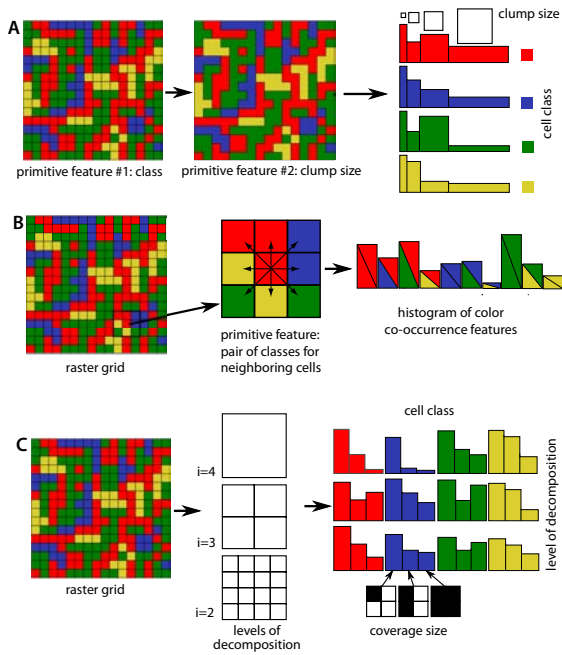


Figure 3: Three of the signature methods implemented in GeoPAT: (A) crossproduct, (B) co-occurrence, (C) decomposition.

as different colors) and clump size ($N_2=4$ categories depicted as increasing size squares). In this example the crossproduct histogram has 16 bins, the value of each bin is a percentage of cells having specified cell class and specified clump size. Crossproduct signature is designed to be effective in encapsulating spatial structures with clear geometric quality (having relatively low complexity); an example of a dataset with such a structure is the land cover raster. In our pattern-based analysis of the NLCD (Jasiewicz et al., 2013; Stepinski et al., 2014; Netzels and Stepinski, 2015) we used the crossproduct signature with two features (cell class with 16 categories and clump size with 14 categories) paired with the Jensen-Shannon divergence distance function (see section 3.2).

co-occurrence – This method uses a color co-occurrence histogram (Barnsley and Barr, 1996; Chang and Krumm, 1999), a variant of the Gray-Level Co-occurrence Matrix (GLCM) originally introduced by Haralick et al. (1973) to characterize texture in grayscale images. In GeoPAT, color is replaced by cell class and a single cell separation of one pixel is used to calculate a co-occurrence histogram. This results in a single primitive feature - a pair of classes assigned to two neighboring cells; eight-connectivity is assumed for establishing the existence of a neighborhood rela-

404 tionship between the two cells. Thus, eight features are
 405 calculated for each cell, but their total number is halved
 406 as the same feature is generated twice by the pairs of
 407 neighboring cells. For a scene with k cell classes, the
 408 co-occurrence histogram has $(k^2 + k)/2$ bins, k of them
 409 correspond to same-class pairs, which measure the com-
 410 position of the classes in the scene, and $(k^2 - k)/2$ bins
 411 correspond to different-class pairs, which measure the
 412 configuration of the classes in the scene. Fig. 3B shows
 413 schematically the construction of a co-occurrence histo-
 414 gram for a scene with $k=4$ classes resulting in a histo-
 415 gram with ten bins. The co-occurrence signature is de-
 416 signed to be effective in encapsulating spatial structures
 417 exhibiting high complexity patterns like the ones result-
 418 ing from a geomorphons-based classification of a DEM
 419 (see Fig. 1). In our pattern-based analysis of DEM data
 420 classified to $k=10$ landform classes (Jasiewicz et al.,
 421 2014) we used the co-occurrence signature with 55 bins
 422 paired with the Wave Hedges distance function (see sec-
 423 tion 3.2).

decomposition – This signature method is inspired
 424 by the work of Remmel and Csillag (2006) to describe
 425 a scene using a set of sub-scenes having a hierarchy of
 426 sizes. For the decomposition method to work best the
 427 scene should be a square having a linear size of 2^D cells.
 428 The scene with k cell classes is scanned without overlap
 429 by a series of square moving windows with sizes $w = 2^i$
 430 cells where $i = 2, \dots, D$ are the decomposition levels.
 431 The size of the maximum scanning window, 2^D , is the
 432 size of the scene. At the smallest decomposition level
 433 $i = 2$ a scene is scanned by a window having a size
 434 of $4 \times 4 = 16$ cells. At each scanning position the per-
 435 centages, p_1, \dots, p_k , of the window's area occupied by
 436 cells having classes $1, \dots, k$, respectively are recorded
 437 and a window area is assigned a list of k tags (one for
 438 each class) representing those percentages. These tags
 439 are classified into one of three categories, 1 if the per-
 440 centage is below $1/4$, 2 if it is between $1/4$ and $1/2$, and 3
 441 if it is above $1/2$. Tallying all tags results in a histogram
 442 with $3 \times k$ bins (three bins for each class).

443 For example, for a scene having a size of 16×16
 444 cells and $k=4$ classes (see Fig. 3C) the number of tags
 445 for decomposition level $i = 2$ is $16 \times 4 = 64$ (number of
 446 sub-windows \times number of classes). These tags are his-
 447 togrammed into 12 bins (number of classes \times number of
 448 tag categories). If, for example, the entire scene is occu-
 449 pied by only one class (say, red), eight bins are equal to
 450 0 and 4 bins (red-3, blue-1, green-1, and yellow-1) have
 451 16 tags each. In this method tags are the primitive fea-
 452 tures. Repeating the same procedure for remaining de-
 453 composition levels results in $D - 1$ histograms each hav-
 454 ing $3 \times k$ bins. All these histograms can be concatenated

456 into a histogram of length equal to $3 \times k \times (D-1)$. Fig. 3C 502
 457 shows schematically a construction of the decomposi- 503
 458 tion histogram for a scene of size 16×16 cells and $k=4$ 504
 459 classes. The size of the scenes dictates the maximum 505
 460 level of decomposition $D=4$ and the histogram length 506
 461 equal to 36. The decomposition signature is designed to
 462 be effective for patterns of all levels of complexity, how-
 463 ever, we have not yet accumulated sufficient experience
 464 working with this signature to offer definitive advice on
 465 the types of datasets to which it can be best applied.

466 3.2. Distance functions

467 *Distance*, which assesses the degree of dissimilarity 510
 468 between two scenes, is the opposite of similarity. The 511
 469 input to all distance functions implemented in GeoPAT 512
 470 is a pair of normalized (the sum of all bins adds to 513
 471 1) signature histograms P and Q and the output is a 514
 472 real number assessing the dissimilarity (distance) be- 515
 473 tween those histograms. When the value of distance 516
 474 function is equal to zero identical histograms are in- 517
 475 dicated, and thus scenes have identical or very similar 518
 476 patterns, whereas large values of the distance function 519
 477 indicate very different histograms and scenes having 520
 478 significantly different patterns. Note that all histogram 521
 479 distance measures are heuristic and no single measure 522
 480 will work well with all signatures. Of over 40 possi- 523
 481 ble histogram distance measures (Cha, 2007) GeoPAT 524
 482 implements the three methods described below which 525
 483 work well with signatures described in the previous sub- 526
 484 section. All three measures have a range of possible 527
 485 values limited to an interval between 0 and 1. 528

486 **Jensen-Shannon divergence** – This measure (Lin, 529
 487 1991) expresses the informational distance between two 530
 488 histograms P and Q by calculating a deviation between 531
 489 the Shannon entropy of the mixture of the two his-
 490 tograms $(P+Q)/2$ (the second term in eq.(1) below) and
 491 the mean of their individual entropies (the second term
 492 in eq.(1)). The value of the Jensen-Shannon divergence
 493 is given by the following formula,

$$d_{JS} = \sqrt{\sum_{i=1}^d \left[\frac{P_i \log_2 P_i + Q_i \log_2 Q_i}{2} - \left(\frac{P_i + Q_i}{2} \right) \log_2 \left(\frac{P_i + Q_i}{2} \right) \right]} \quad (1)$$

494 where d is the number of bins (the same for both his- 532
 495 tograms) and P_i and Q_i are the values of i th bin in the 533
 496 two histograms. We have found Jensen-Shannon diver- 534
 497 gence works well (yields dissimilarity values in agree- 535
 498 ment with human visual perception) for comparison of 536
 499 land cover patterns as encapsulated by crossproduct sig- 537
 500 natures (Jasiewicz et al., 2013; Stepinski et al., 2014; 538
 501 Netzel and Stepinski, 2015).

Wave Hedges – This measure is designed to work
 with co-occurrence signature histograms that tend to be
 dominated by bins corresponding to adjacent cells hav-
 ing the same class. The value of the Wave Hedges dis-
 tance is given by the following formula (Cha, 2007),

$$d_{WH} = \sum_{i=0}^d e_i \frac{|P_i - Q_i|}{\max(P_i, Q_i)} \quad (2)$$

where d is the maximum number of possible bins and
 $e_i = 1$ if $\max(P_i, Q_i) > 0$ or $e_i = 0$ otherwise. In
 other words, only pattern features present in at least
 one of the two scenes contribute to the value of the
 distance. In Wave Hedges distance formula all present
 features contribute to the overall value of distance with
 the same weights regardless of feature abundance in
 the scenes. In the case of the co-occurrence his-
 togram (Fig. 3B) this means that composition-related
 features and configuration-related features contribute
 equally to the distance value despite the heavy domi-
 nance of composition-related features in histograms
 stemming from all realistic scenes. This makes the
 Wave Hedges distance particularly suitable for compar-
 ison of terrain scenes as encapsulated by co-occurrence
 signatures (Jasiewicz et al., 2014)

Jaccard– This measure is an extension of the Jaccard
 similarity coefficient (Jaccard, 1908), originally devel-
 oped to assess a similarity between two sets, but used
 here for assessing the dissimilarity between two his-
 tograms. The value of the Jaccard distance is given by
 the following formula (Cha, 2007),

$$d_J = 1 - \frac{\sum_{i=1}^d P_i Q_i}{\sum_{i=1}^d P_i^2 + \sum_{i=1}^d Q_i^2 - \sum_{i=1}^d P_i Q_i} \quad (3)$$

We have found that the Jaccard distance works well for
 comparison of land cover patterns as encapsulated by
 decomposition signatures.

532 4. Core modules

The seven core modules in the GeoPAT toolbox pro-
 vide the infrastructure for pattern-based analysis of spa-
 tial datasets. There are two types of core modules: sig-
 nature extraction modules and geoprocessing modules.
 Fig. (4) illustrates different possible pipelines of data
 processing using these modules.

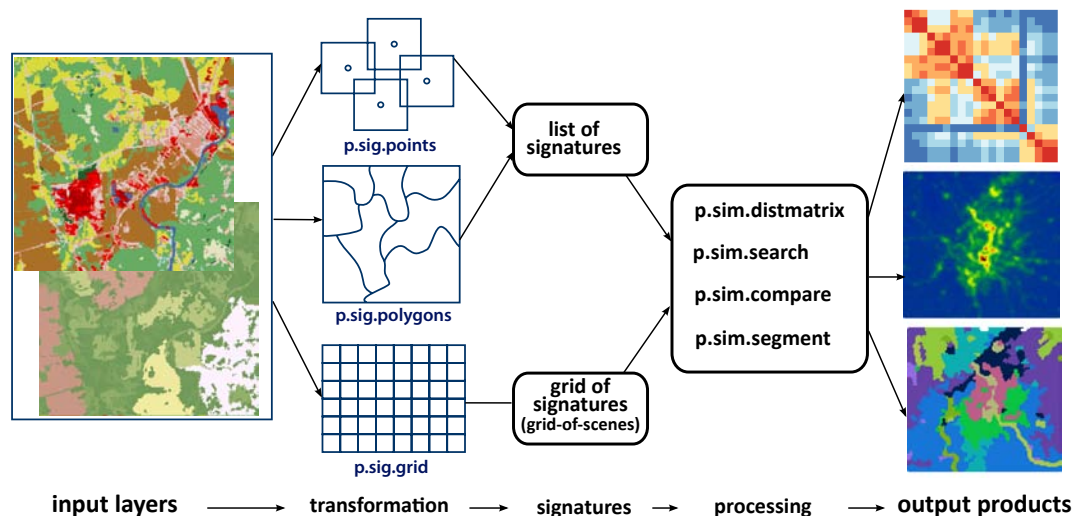


Figure 4: Different data processing pipelines possible using modules from the GeoPAT toolbox. Input consists of raster layers and output consists of raster layers or text tables.

4.1. Signature extraction modules

The input to all three signature extraction modules is a set of categorical rasters containing all information layers needed to construct scenes signatures. The output is a set of signatures; each module outputs these signatures in a different data structure depending on its definition of a scene or a set of scenes. Fig. 5 shows three possible scenarios for scene definition which are addressed by the three modules.

p.sig.points. This module extracts signatures for a collection of individual scenes having a square geometry (see Fig. 5A). The user provides the coordinates of the center of each scene (point file) and the size of the scene. The module outputs a list of scene-labeled signatures. Note that p.sig.points can be used to extract a signature for the entire region if needed. There are several typical uses for this module. Examples include generating a query scene to be compared with a grid-of-scenes for the search task, comparison of several scenes in a single raster (like in a comparison of different cities on the basis of their patterns of land cover classes), and comparison of two different co-registered rasters (like in comparison of a natural scene with a scene resulting from a computer simulation aimed at recreating a pattern observed in the natural scene).

p.sig.polygons. This module extracts signatures for a collection of individual scenes having a polygonal geometry (see Fig. 5B). The user provides as input a categorical raster layer which defines the division of a region into polygonal scenes and the module outputs a

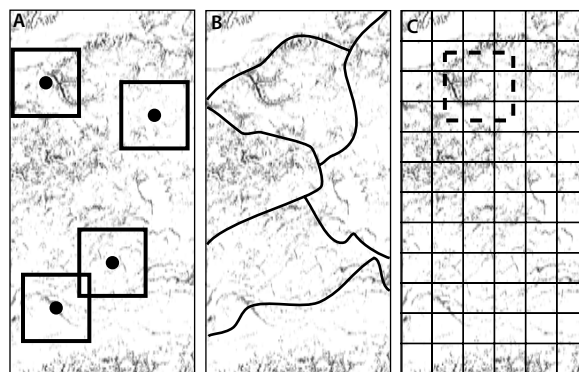


Figure 5: Methods of scene definition: A – scenes defined by points; B – scenes defined by polygons; C – scenes defined by a grid.

list of polygon-labeled signatures. A typical use for p.sig.polygons is for comparing irregular scenes resulting from a segmentation of the region (using the segmentation module p.sim.segment, see the next subsection).

p.sig.grid. This module extracts a grid-of-scenes (see Fig. 5C) – a grid of the same spatial extent as the region defined by the input data but having larger cells (s-cells). Each s-cell has only one attribute - a signature of the scene centered on it. The module outputs a header file containing the topology of the grid-of-scenes and a binary file containing signatures ordered row by row. The grid-of-scenes is an input to three geoprocessing modules: p.sim.search, p.sim.segment and p.sim.compare.

4.2. Geoprocessing modules

The four geoprocessing modules use scene signatures and distance functions to perform four popular analysis tasks including: comparison of individual scenes (`p.sim.distmatrix`), comparison between a single scene (a query) and grid-of-scenes (`p.sim.search`), comparison between two grids-of-scenes (`p.sim.compare`), and a segmentation of a grid-of-scenes (`p.sim.segment`). Fig. (6) illustrates the tasks performed by these modules.

p.sim.distmatrix. This module computes a distance matrix between a collection of scenes. It uses signatures output by `p.sig.points` or `p.sig.polygons` modules and an appropriate distance function from the library. The resultant distance matrix is typically used as an input for scene clustering (Fig. 6A) which results in discovering structures in the data without guidance from an analyst. Clustering itself is not implemented in GeoPAT nor is it implemented directly in GRASS; we recommend using the hierarchical clustering algorithm implemented in R (R Core Team, 2013) as GRASS is designed to work together with R (Bivand, 2000). The distance matrix generated by `p.sim.distmatrix` is the only required input to the hierarchical clustering algorithm. An example of `p.sim.distmatrix` usage would be the clustering of a collection of cities on the basis of the patterns of land cover classes within their boundaries.

p.sim.search. This module performs a query-by-pattern-similarity (QBPS). The input is a query scene (or a list of query scenes) output by `p.sig.points` or `p.sig.polygons` modules and a grid-of-scenes (database to be queried) output by the `p.sig.grid` module. The signature of each query and the signatures in a grid-of-scenes must have the same structure. The module compares query/queries with the database in a scene-by-scene fashion and outputs a layer(s) having the same topology as the grid-of-scenes but containing values of similarity between a query and each scene in the grid-of-scenes. The results of QBPS can be visualized as a similarity map (Fig. 6B).

QBPS provides a knowledge discovery tool that is qualitatively different from retrieval of top-matches-to-a-query (Câmara et al., 1996; Datcu et al., 2002; Kopferski et al., 2002; Aksoy et al., 2005; Barb and Shyu, 2010) – a standard approach to searching for similar scenes. QBPS is a GIS tool inasmuch as it performs spatial processing resulting in a map that shows geographical distribution of degree of similarity to the query scene. Such map provides much more information than a non-spatial list of top matches to a query. By utilizing spatial organization it simultaneously shows similarity relations between the query and *all* scenes in the

database. Thus it allows an analyst to concentrate on revealed geospatial phenomena rather than on similarity between specific scenes. QBPS has been used to query the NLCD 2006 dataset (Jasiewicz and Stepinski, 2013a; Stepinski et al., 2014) for similarity between land cover scenes and to query topography of the country of Poland (Jasiewicz et al., 2014) for similarity between landscapes. Note that QBPS can be used as an element of supervised classification of a region into different pattern types such as, for example, urban structures, landscape types, or physiographic units (see the next section).

p.sim.compare. This module compares two grids-of-scenes in a scene-by-scene fashion. The grids must have the same topologies and the output of `p.sim.compare` is a raster layer having the same topology as the inputs but containing values of similarity between corresponding pairs of scenes (Fig. (6C)). This is equivalent to the GIS *overlay* function, but is performed on signature attributes. Applying `p.sim.compare` to two land cover datasets of the same region but pertaining to two different time steps enables pattern-based change detection (Netzel and Stepinski, 2015). Unlike traditional land cover change detection which tracks cell-by-cell transitions of land cover categories, pattern-based change detection assesses change in local patterns of land cover; it is especially useful for continental-scale or global-scale assessments of land cover change. Another application of `p.sim.compare` is for comparison of two layers created using different parameters. For example, the module can be used for comparison of two classifications of a DEM using the same method and the same target landscape classes but different extraction parameters (Jasiewicz and Stepinski, 2013b). Such a comparison allows a user to see what landscape types are most sensitive to the values of free parameters in their mapping algorithm.

p.sim.segment. This module segments a grid-of-scenes into regions of uniform patterns (Niesterowicz and Stepinski, 2013) (Fig. 6D) in the same fashion as traditional segmentation algorithms segment image (or other rasters) into segments of uniform color and texture. The difference is that `p.sim.segment` segments a grid-of-scenes rather than an ordinary grid, and uses scene signature as attribute rather than color or image texture to decide how to delineate the segments. The module uses a variant of the region growing algorithm (Zucker, 1976; Câmara et al., 1996; Li and Narayanan, 2004; Blaschke, 2010). It has two free parameters: a similarity threshold that defines the minimum similarity between two scenes to be treated as “similar” and, optionally, a minimum number of scenes to constitute a

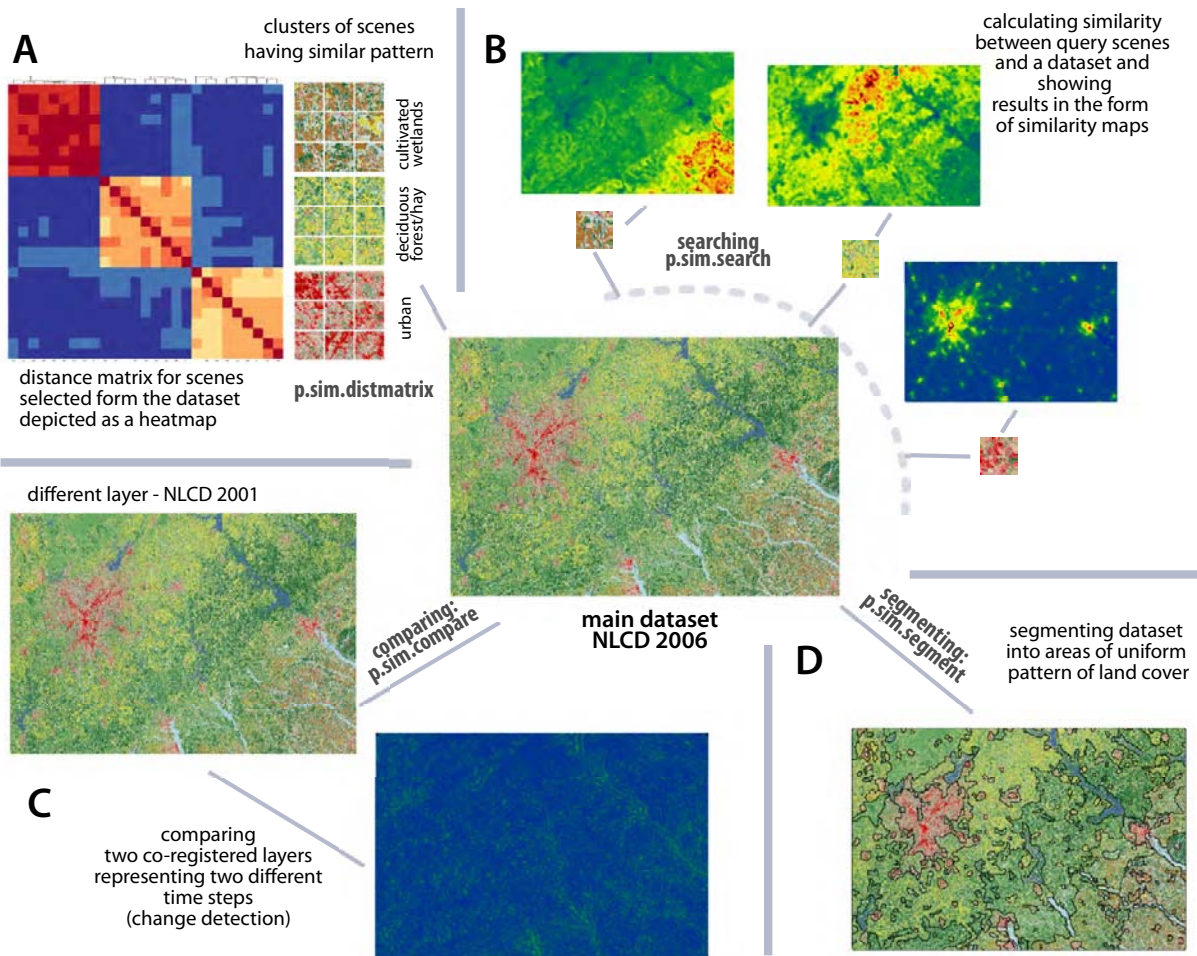


Figure 6: Illustration of data processing, see text for details

687 separate segment.

688 The segmentation provided by the p.sim.segment
 689 module is typically used as an intermediate step in regionalization of the dataset. The goal of regionalization
 690 is to generalize and thus simplify spatial representation
 691 of data so it is more meaningful and easier to analyze.
 692 Examples of regionalization include delineation of landscape types (Niesterowicz and Stepinski, 2013) or delineation
 693 of physiographic units (Jasiewicz et al., 2014).
 694 Regionalization is achieved by clustering segments output by p.sim.segment into a number of distinct pattern type classes.
 695
 696
 697
 698

699 5. Case study

700 To demonstrate GeoPAT's capabilities in application
 701 to a specific dataset we use a region extracted from

702 the NLCD 2006 referred to as "Atlanta." Atlanta covers a region in the northern part of the U.S. state of Georgia (see Fig. 8A) that includes the city of Atlanta.
 703 The grid has a size of 11300 x 7500 cells and a resolution of 30m/cell. It is a categorical grid with 16 land cover classes. The data can be downloaded from <http://sil.uc.edu>. NLCD is the result of classification of Landsat images, so this example pertains to image data. Working with the NLCD avoids performing the pre-processing step of pixel-based classification, which is not a part of GeoPAT toolbox.
 704
 705
 706
 707
 708
 709
 710
 711
 712

713 The purpose of this case study is to perform unsupervised and supervised regionalizations of the Atlanta site into landscape types (characteristic patterns of land cover) using GeoPAT modules and R. Two different modes of machine learning (unsupervised and supervised) are demonstrated to show the range of tasks that
 714
 715
 716
 717
 718

719 can be achieved using GeoPAT. The schema of the two 750
 720 procedures are shown in Fig. 7. These procedures can 751
 721 be run as a single routine due to the full integration be- 752
 722 tween GRASS 7 and R (Bivand et al., 2008). These 753
 723 routines are available as supplementary material to this 754
 724 paper (<http://sil.uc.edu>). Based on our earlier experi- 755
 725 ence in working with the NLCD (Jasiewicz and Stepin- 756
 726 ski, 2013a; Stepinski et al., 2014; Netzel and Stepin- 757
 727 ski, 2015) we use the crossproduct signature for scene 758
 728 representation and the Jensen-Shannon divergence as a 759
 729 distance function to measure dissimilarity between the 760
 730 scenes. The first step in both procedures is to gener- 761
 731 ate the grid-of-scenes using the p.sig.grid module. We 762
 732 have selected the grid-of-scenes to have a s-cell equal to 763
 733 900m (30 times the size of cell in the Atlanta grid) and 764
 734 we define scenes as square regions having the size of 4.5 765
 735 km×4.5 km. Thus, the scenes overlap significantly.

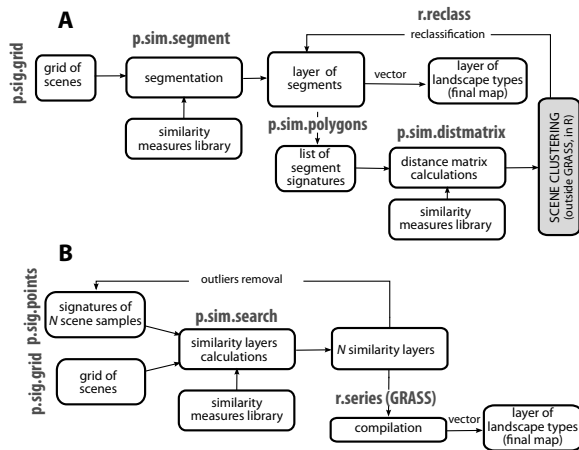


Figure 7: Schemes of processing pipelines for case study calculations: A) unsupervised regionalization; B) supervised regionalization

5.1. Unsupervised regionalization

736 We start by running the p.sim.segment module with 737
 738 the following parameters: similarity threshold=0.75, 739
 739 minimum number of scenes=10. This yielded 434 indi- 740
 740 vidual segments. Signatures of the segments were cal- 741
 741 culated using the p.sim.polygons module. Using these 742
 742 signatures and the Jensen-Shannon divergence we clus- 743
 743 ter the 434 segments into seven “landscape types” using 744
 744 hierarchical clustering algorithm with “Ward” linkage 745
 745 (available in R “stats” package). The choice of a num- 746
 746 ber of clusters is arbitrary as is always the case in hi- 747
 747 erarchical clustering. The result is a map of landscape 748
 748 types in the Atlanta site (Fig. 8B). Note that these lan- 749
 749 dscape types emerged from the data and we assign labels

750 to them (see caption to Fig. 8) only a posteriori, after 751
 751 reviewing the results of hierarchical clustering.

752 In principle, an unsupervised classification could 753
 753 be performed without the intermediate segmentation 754
 754 step. The signatures calculated by the p.sig.grid mod- 755
 755 ule could be clustered to yield landscape types, however 756
 756 this would involve clustering over 90,000 scenes and 757
 757 the results would exhibit salt and pepper noise. The pro- 758
 758 cedure demonstrated here employs the concept of object- 759
 759 based analysis – first segment then classify but it is the 760
 760 grid-of-scenes rather than an image which is segmented.

5.2. Supervised regionalization

761 To begin, we selected seven individual scenes 762
 762 (shown as white squares in Fig. 8C) as landscape type 763
 763 archetypes or samples of the seven landscape types we 764
 764 have chosen for mapping; 1 – urban, 2 – wet crop- 765
 765 lands, 3 – pasture-dominated, 4 – deciduous forest- 766
 766 dominated, 5 – evergreen forest-dominated, 6 – wet- 767
 767 lands and surroundings, 7 – waters and surroundings. 768
 768 Signatures for the seven sample scenes were calculated 769
 769 using the r.sig.points module. These signatures were 770
 770 used as queries over the Atlanta grid-of-scenes using 771
 771 the p.sim.search module resulting in seven similarity 772
 772 maps, one for each query. The maps were overlaid (us- 773
 773 ing GRASS capabilities) and each s-cell was assigned a 774
 774 landscape type label corresponding to the largest value 775
 775 of similarity resulting in a single map of landscape types 776
 776 (Fig. 8D).

777 It is interesting to observe similarities and differences 778
 778 between the two maps. Differences are expected be- 779
 779 cause the two maps were obtained following very dif- 780
 780 ferent principles. The unsupervised map on Fig. 8B re- 781
 781 flects the natural grouping of landscapes subject to a re- 782
 782 striction on the total number (in our case – seven) of 783
 783 the groups. The supervised map on Fig. 8D reflects the 784
 784 preferences of an analyst who has selected a priori spe- 785
 785 cific landscape types to be mapped. The similarities be- 786
 786 tween the two maps stem from the fact that both of them 787
 787 reflects the same physical reality.

788 The two processing schemes demonstrated here 789
 789 closely resemble standard schemes for unsupervised 790
 790 and supervised image classifications. This is intentional 791
 791 as GeoPAT is designed to work like a standard GIS sys- 792
 792 tem but with scenes rather than pixels. However, the 793
 793 output of GeoPAT (see, for example, Figs. 8B and D) 794
 794 is fundamentally different from the output of an image 795
 795 classification algorithm as it yields pattern-based gen- 796
 796 eralization of pixel-based classification.

797 The unsupervised example illustrates the concept of 798
 798 generalization that has been previously applied (using 799
 799 different methods) to land cover datasets in the context 800

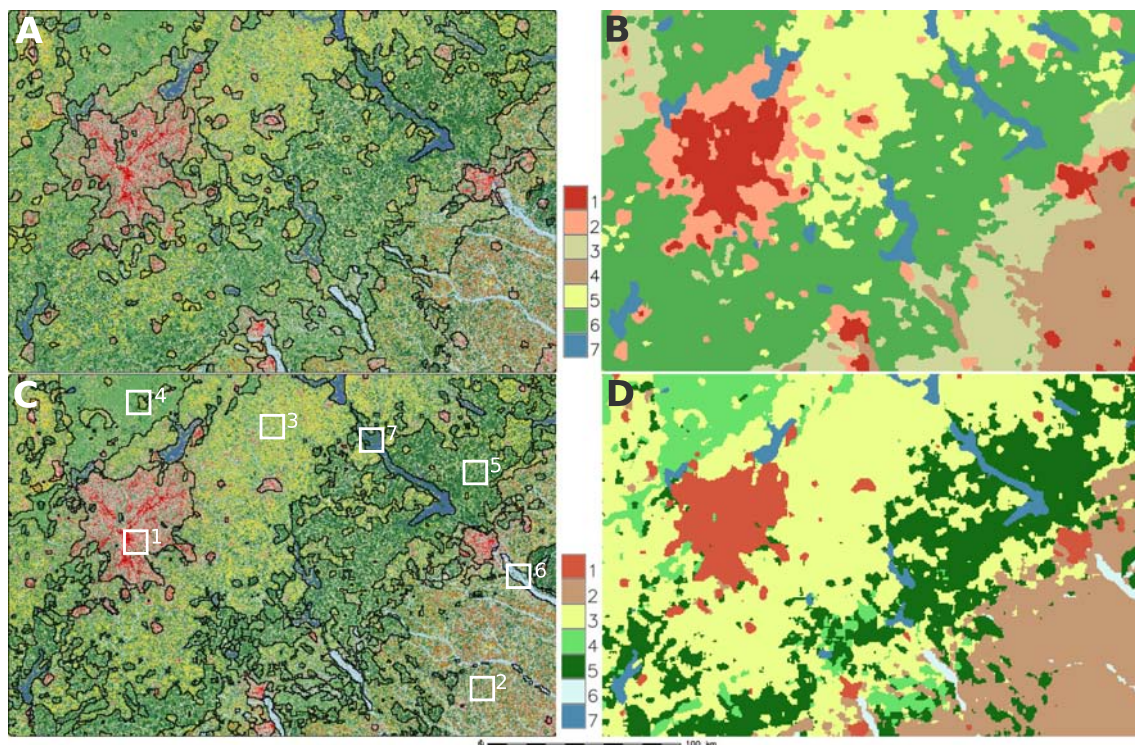


Figure 8: Results of the case study. Boundaries of delineated landscape types obtained using unsupervised (A) and supervised (C) regionalization procedures superimposed on the NLCD map (see <http://www.mrlc.gov/> for the land cover legend). Regionalization maps resulting from unsupervised (B) and supervised (D) classification, respectively. Classes of unsupervised regionalization are: 1 – urban, 2 – suburban, 3 – pasture-forest mixture, 4 – wet croplands, 5 – pasture-dominated, 6 – forests dominated, and 7 – waters and surroundings. Classes of supervised classification are: 1 – urban, 2 – wet croplands 3 – pasture-dominated, 4 – deciduous forest-dominated, 5 – evergreen forest-dominated, 6 – wetlands and surroundings, and 7 – waters and surroundings. Scene samples utilized in supervised approach are marked by white squares.

801 of landscape ecology (Long et al., 2010; Cardille and
 802 Lambois, 2009) and the supervised example illustrates a
 803 generalization concept that has been previously applied
 804 (using different methods) to RS images in the context of
 805 mapping urban landscapes (Moller-Jensen et al., 2005;
 806 Vatsavai, 2013a; Graesser et al., 2012). With GeoPAT
 807 these kinds of generalizations can be performed with
 808 relative ease and on much larger datasets by taking advan-
 809 tage of GRASS' ability to handle very large datasets.

810 6. Performance

811 GeoPAT has been optimized to work efficiently with
 812 big data where it is most effective as a knowledge dis-
 813 covery tool. It can be applied to giga-cell rasters when
 814 running on either servers or workstations. Table 1
 815 lists execution times for GeoPAT modules as applied to
 816 several large datasets. This gives a rough idea about
 817 GeoPAT's level of performance. All calculations were
 818 run on a double-CPU XEON machine (8 cores each)
 819 with 20 GB of RAM running Linux. Three datasets

820 were used: (1) The POLAND dataset ($24,000 \times 27,000$
 821 cells) is a 10-classes map of landform elements calcu-
 822 lated from a 30m/cell DEM (Jasiewicz and Stepin-
 823 ski, 2013b); (2) The CHINA dataset ($84,000 \times 64,000$
 824 cells) is a 10-class map of landform elements calcu-
 825 lated from a 90m/cell DEM (SRTM); (3) The USA dataset
 826 ($164,000 \times 104,000$ cells) is a 16-class, 30m/cell map
 827 of land cover/land use (NLCD 2006) covering the entire
 828 conterminous United States.

829 Table 1 is divided into two parts, part A pertains to
 830 the performance of the signature extraction modules and
 831 part B pertains to the performance of the geoprocessing
 832 modules. In general, the signature extraction modules
 833 are significantly more computationally expensive than
 834 geoprocessing modules. However, in a typical appli-
 835 cation signature extraction needs to be performed only
 836 once, whereas geoprocessing computation may require
 837 multiple runs in order to get a satisfactory result.

838 Note that the p.sim.search module is fast enough to
 839 enable real-time search. Indeed, this module provides
 840 a computational engine for our two GeoWeb pattern

module	input	output	processing time
A. Signature extraction modules			
p.sig.points with co-occurrence function	single scene 300×300 cells with 10 categories	single signature with 55 bins	0s 73ms
p.sig.grid with co-occurrence function	POLAND	800×900 grid-of-scenes	2h 26m
p.sig.grid with co-occurrence function	CHINA	1680×1280 grid-of-scenes	20h 39m
p.sig.grid with crossproduct function parallelized with 10 threats	USA	1640×1040 grid-of-scenes	4h 46m
B. Geoprocessing modules			
p.sim.distmatrix with Wave-Hedges function	1084 histograms with 136 bins each	1084×1084 distance matrix	48s
p.sim.search , 55-bins histograms, Wave-Hedges function	64 queries, POLAND grid-of-scenes	64 800×900 similarity layers	5s 33ms
p.sim.search , 192-bins histograms, Jensen-Shannon function	1 query, USA grid-of-scenes	one 1640×1040 similarity layer	6s 11ms
p.sim.compare , 55-bins histograms, Wave-Hedges function	two POLAND grids-of-scenes	one 800×900 similarity layer	0s 97ms
p.sim.segment , 55-bins histograms, Wave-Hedges function	one POLAND grid-of-scenes	one 800×900 layer containing segments	1s 28ms

Table 1: Examples of computation times for different modules of GeoPAT and using different datasets.

841 search applications, LandEx-USA – which enables dis- 862
842 discovery of similar land cover patterns over the extent of 863
843 the United States and TerraEx-PL – which enables dis- 864
844 discovery of similar landscapes over the extent of the coun- 865
845 try of Poland. These online applications are available at 866
846 <http://sil.uc.edu/>. The real-time response to a query in 867
847 these applications is achieved by pre-calculating grids- 868
848 of-scenes and storing them in the RAM. Once a query 869
849 is submitted by a user the application runs p.sim.search 870
850 and returns the generated similarity map. 871

851 7. Discussion and conclusions

852 GeoPAT is a toolbox which implements our method 872
853 (Jasiewicz and Stepinski, 2013b; Stepinski et al., 2014; 873
854 Jasiewicz et al., 2014; Netzel and Stepinski, 2015) for 874
855 pattern-based information retrieval from images and 875
856 other rasters. The recent interest in such methods stems 876
857 from the need to consider spatial patches larger than 877
858 a pixel to adequately reflect local content. To the 878
859 best of our knowledge, the only current methodological 879
860 (but not software) alternative to GeoPAT is the Com- 880
861 plex Object-Based Image Analysis (COBIA) (Vatsavai, 881

862 2013a,b) which relies on the MIL concept. At the con- 882
863 ceptual level GeoPAT and COBIA are similar inasmuch 883
864 as both use scenes, both describe a scene as a Probabil- 884
865 ity Distribution Function (PDF) of image features, and 885
866 both use measures of similarity between PDFs to asses 886
867 a degree of similarity between the scenes. The differences 887
868 are in an the implementation of this concepts. COBIA 888
869 works directly with images and assumes that the PDF of 889
870 image features can be modeled by a multi-variate Gaus- 890
871 sian. GeoPAT works with previously classified images 891
872 and models PDFs as histograms. Both of these meth- 892
873 ods differ from earlier approaches (Moller-Jensen et al. 893
874 (2005) or Lucieer and Stein (2005)) to improve clas- 894
875 sification and segmentation of images by incorporating 895
876 texture descriptors as additional features in pixel-based 896
877 (single-instance learning) algorithms. 897

878 The decision to design GeoPAT to operate on cat- 898
879 egorical rasters was dictated by two considerations: 899
880 (a) the flexibility of input data, so different modalities 900
881 of raster data could be analyzed by GeoPAT once 901
882 classified, and (b) the use of categorical rasters in- 902
883 creases the performance of an algorithm and make 903
884 it easier to work with large datasets. Note that 904
885 many large datasets of interest are already avail- 905

886 able in the categorized (land cover) form. These 938
887 include 30m U.S.-wide NLCD, 30m global GLC30 939
888 (<http://www.globallandcover.com/>), 300m global 940
889 GlobCover (<http://due.esrin.esa.int/globcover/>), and 941
890 100m Europe-wide CORINE. In addition, topographic 942
891 datasets (DEMs) can be categorized with relative ease 943
892 by using robust techniques, such as the geomorphons 944
893 method (Jasiewicz and Stepinski, 2013b) for which an 945
894 open source code as well as an online application are 946
895 available at <http://sil.uc.edu/>. Additionally, RGB VHR 947
896 images can be categorized by clustering colors with 948
897 Euclidean distance in CIE-Lab color space (Rubner 949
898 et al., 2000). 950

899 Given a dataset, what combination of signature/distance 951
900 function are most appropriate? There is no quantitative 952
901 means to assist in answering this question directly. 953
902 Each combination is based on different set of 954
903 heuristics and will yield different values of similarities 955
904 for the same pairs of scenes. These values can be compared 956
905 (Stepinski and Cohen, 2014) by using their corresponding 957
906 percentiles (calculated from the empirical cumulative 958
907 distribution functions of the sets of similarities between 959
908 all scenes in the dataset), but it is not possible to determine 960
909 objectively which similarity value is “better” as this invokes 961
910 subjective human perception of similarity between spatial 962
911 patterns. 963

912 However, in applications of GeoPAT to classification 964
913 tasks, the utility of a particular signature/distance function 965
914 combination can be assessed indirectly by assessing quality 966
915 of the classification. When using the GeoPAT for supervised 967
916 learning, the quality of a classifier can be assessed by applying 968
917 the classifier to a test set of pre-labeled scenes and calculating 969
918 the standard metrics of performance. The results will depend on 970
919 the selection of a particular signature/distance function combination. 971
920 Note that when a region under consideration has clearly visible 972
921 divisions between different pattern types, like in the case of 973
922 the DEM shown in Fig. 1 or the image used by Vatsavai (2013a), 974
923 all possible classifiers will perform very well as the problem 975
924 presents little challenge to an algorithm. On the other hand, 976
925 regions on which an analyst would have difficulty in delineating 977
926 the boundaries between different types of patterns, like in the 978
927 case of Atlanta region, present a challenge. First, it is difficult- 979
928 to-impossible to manually delineate different pattern types in 980
929 such regions (to construct a test set), second, different classifiers 981
930 (different combinations of signature/similarity function) will 982
931 yield different results without an objective means to assess their 983
932 quality due to lack of a reliable test set. However, in our opinion, 984
933 such “difficult” regions are where GeoPAT is most useful as it 985
934 retrieves information which could be difficult to retrieve 986
935 987
936 988
937 989

by any other means. To make this point clearer consider 938
the Atlanta region as shown in the middle of Fig. 6. It 939
would be very difficult for an analyst to manually delineate 940
pattern types in this region and no two manual delineations 941
would be the same, but GeoPAT delineations (Fig. 8) are 942
objective, repeatable, and they make perfect sense to an analyst 943
once the computed boundaries are superimposed on the land cover 944
map (Fig. 8 A and C). 945

Finally, GeoPAT is an actively developed software and we 946
expect users to contribute to it by adding to the shared library 947
of functions. GeoPAT is available for download from <http://sil.uc.edu/>. 948
Currently, it has been tested on and made available for the Linux 949
operating system; it requires the development version of GRASS 950
7. 951

953 Acknowledgments

The authors wish to thank J. Niesterowicz and A. 954
Dmowska for helpful comments and discussions. The work was 955
supported by the National Science Centre (NCN) under Grant 956
DEC-2012/07/B/ST6/01206, and by the National Science Foundation 957
(NSF) under Grant BCS- 621 1147702, and by the University of Cincinnati 958
Space Exploration Institute. 959
960

961 References

- 962 Aksoy, S., Koperski, K., Tusk, C., Marchisio, G., Tilton, J., Mar. 2005.
963 Learning bayesian classifiers for scene classification with a visual
964 grammar. *IEEE Trans. Geosci. Remote Sens.* 43 (3), 581–589.
965 Barb, A., Shyu, C., 2010. Visual-Semantic Modeling in Content-
966 Based Geospatial Information Retrieval Using Associative Mining
967 Techniques. *IEEE Geosci. Remote Sens. Lett.* 7 (1), 38–42.
968 Barnsley, M., Barr, S., 1996. Inferring urban land use from satellite
969 sensor images using kernel-based spatial reclassification. *Photogrammetric engineering and remote sensing* 62 (8), 949–958.
970 Bivand, R. S., 2000. Using the r statistical data analysis language
971 on grass 5.0 gis database files. *Computers & Geosciences* 26 (9),
972 1043–1052.
973 Bivand, R. S., Pebesma, E. J., Gómez-Rubio, V., 2008. *Applied spatial
974 data analysis with R*. Vol. 747248717. Springer.
975 Blaschke, T., Jan. 2010. Object based image analysis for remote sensing. *ISPRS J. Photogramm. Remote Sens.* 65 (1), 2–16.
976 Cain, D., Riitters, K., 1997. A multi-scale analysis of landscape statistics. *Landscape Ecol.* (630), 199–212.
977 Câmara, G., Souza, R. C. M., Freitas, U. M., Garrido, J., May 1996.
978 Spring: Integrating remote sensing and gis by object-oriented data
979 modelling. *Comput. Graph.* 20 (3), 395–403.
980 Cardille, J. A., Lambois, M., 2009. From the redwood forest to the
981 Gulf Stream waters: human signature nearly ubiquitous in representative
982 US landscapes. *Frontiers in Ecology and the Environment* 8(3), 130–134.
983 Cardille, J. A., White, J. C., Wulder, M. A., Holland, T., 2012. Representative
984 landscapes in the forested area of Canada. *Environmental management* 49(1), 163–173.
985
986
987
988
989

- 990 Cha, S., 2007. Comprehensive survey on distance/similarity measures 1055
991 between probability density functions. *Int. J. Math. Model. Meth-* 1056
992 *ods Appld Sci.* 1(4) (4), 300–307. 1057
- 993 Chang, P., Krumm, J., 1999. Object recognition with color cooccur- 1058
994 ence histograms. In: *In Proceedings of IEEE Conference on Com-* 1059
995 *puter Vision and Pattern Recognition, Fort Collins, CO, June 23-* 1060
996 *25, 1999. IEEE Computer Society Conference.* 1061
- 997 Cushman, S. A., McGarigal, K., Neel, M. C., 2008. Parsimony in 1062
998 landscape metrics: strength, universality, and consistency. *Ecolog-* 1063
999 *ical indicators* 8(5), 691–703. 1064
- 1000 Daschiel, H., Datcu, M., Jan. 2005. Information mining in remote 1065
1001 sensing image archives: system evaluation. *IEEE Trans. Geosci.* 1066
1002 *Remote Sens.* 43 (1), 188–199. 1067
- 1003 Datcu, M., Daschiel, H., Pelizzari, a., Quartulli, M., Galoppo, a., Co- 1068
1004 lapicchioni, a., Pastori, M., Seidel, K., Marchetti, P., D’Elia, S., 1069
1005 Dec. 2003. Information mining in remote sensing image archives: 1070
1006 system concepts. *IEEE Trans. Geosci. Remote Sens.* 41 (12), 1071
1007 2923–2936. 1072
- 1008 Datcu, M., Seidel, K., D’Elia, S., Marchetti, P., 2002. Knowledge- 1073
1009 driven information mining in remote-sensing image archives. *ESA* 1074
1010 *Bull.* 110 (may), 26–33. 1075
- 1011 Datta, R., Joshi, D., Li, J., Wang, J. Z., Apr. 2008. Image Retrieval: 1076
1012 Ideas, Influences, and Trends of the New Age. *ACM Comput. Surv.* 1077
1013 40 (2), 1–60. 1078
- 1014 Dilts, T. E., Yang, J., Weisberg, P. J., 2010. The landscape similarity 1079
1015 toolbox: new tools for optimizing the location of control sites in 1080
1016 experimental studies. *Ecography* 33(6), 1097–1101. 1081
- 1017 Gevers, T., Smeulders, A. W., 2004. Content-based image retrieval: 1082
1018 An overview. In: Kang, G. M. S. B. (Ed.), *Emerging Topics in* 1083
1019 *Computer Vision. Upper Saddle River, NJ: Prentice-Hall, Ch. 8,* 1084
1020 *pp. 333–384.* 1085
- 1021 Graesser, J., Cheriyyad, A., Vatsavai, R., Chandola, V., Long, J., 1086
1022 Bright, E., 2012. Image based characterization of formal and in- 1087
1023 formal neighborhoods in an urban landscape. *IEEE Journal of Se-* 1088
1024 *lected Topics in Applied Earth Observations and Remote Sensing* 1089
1025 5(4), 1164–1176. 1090
- 1026 GRASS Development Team, 2012. *Geographic Resources Analysis* 1091
1027 *Support System (GRASS GIS) Software. Open Source Geospatial* 1092
1028 *Foundation, USA.* 1093
- 1029 Haines-Young, R., Chopping, M., 1996. Quantifying landscape struc- 1094
1030 ture: a review of landscape indices and their application to forested 1095
1031 landscapes. *Progress in Physical Geography* 20(4), 418–445. 1096
- 1032 Haralick, R. M., Shanmugam, K., Dinstein, I., Nov. 1973. Textural 1097
1033 features for image classification. *Syst. Man Cybern. IEEE Trans.* 1098
1034 3 (6), 610–621. 1099
- 1035 Jaccard, P., 1908. *Nouvelles recherches sur la distribu- tion florale.* 1100
1036 *Bull. Soc. Vaudoise Sci. Nat.* 44, 223–270. 1101
- 1037 Jasiewicz, J., Netzel, P., Stepinski, T. F., 2014. Landscapes similarity, 1102
1038 retrieval, and machine mapping of physiographic units. *Geomor-* 1103
1039 *phology* 221, 104–112. 1104
- 1040 Jasiewicz, J., Stepinski, T., Netzel, P., 2013. Content-based landscape 1105
1041 retrieval using geomorphons. In: *Geomorphometry 2013. Nanjing,* 1106
1042 *China, pp. 1–4.* 1107
- 1043 Jasiewicz, J., Stepinski, T. F., 2013a. Example-Based Retrieval 1108
1044 of Alike Land-Cover Scenes From NLCS2006 Database. *IEEE* 1109
1045 *Geosci. Remote Sens. Lett.* 10 (1), 155–159. 1110
- 1046 Jasiewicz, J., Stepinski, T. F., 2013b. Geomorphons-a pattern recog- 1111
1047 nition approach to classification and mapping of landforms. *Geo-* 1112
1048 *morphology* 182, 147–156. 1113
- 1049 Koperski, K., Marchisio, G., Aksoy, S., Tusk, C., 2002. VisiMine: in- 1114
1050 teractive mining in image databases. In: *IEEE Int. Geosci. Remote* 1115
1051 *Sens. Symp. Vol. 3. Ieee, pp. 1810–1812.* 1116
- 1052 Körting, T., Fonseca, L. G., Câmara, G., 2013. GeoDMAGeographic 1117
1053 Data Mining Analyst. *Comput. Geosci.* 57, 133–145. 1118
- 1054 Kupfer, J. A., Gao, P., Guo, D., May 2012. Regionalization of forest 1119
pattern metrics for the continental United States using contiguity
constrained clustering and partitioning. *Ecol. Inform.* 9, 11–18.
- Lang, S., 2008. Object-based image analysis for remote sensing appli-
cations: modeling reality–dealing with complexity. *Object-Based
Image Anal.*, 3–27.
- Lew, M., Sebe, N., Lifi, C., Jain, R., 2006. Content-based multimedia
information retrieval: State of the art and challenges. *ACM Trans.
Multimedia Comput., Commun., Appl.*, 2(1), 1–19.
- Li, J., Narayanan, R., 2004. Integrated spectral and spatial information
mining in remote sensing imagery. *IEEE Trans. Geosci. Remote
Sens.* 42 (3), 673–685.
- Lin, J., 1991. Divergence measures based on the Shannon entropy.
IEEE Transactions on Information Theory 31(1), 145–151.
- Long, J., Nelson, T., Wulder, M., Jul. 2010. Regionalization of land-
scape pattern indices using multivariate cluster analysis. *Environ.
Manage.* 46 (1), 134–42.
- Lu, D., Weng, Q., 2007. A survey of image classification methods and
techniques for improving classification performance. *International
Journal of Remote Sensing* 28(5), 823–870.
- Lucieer, A., Stein, A., 2005. Texture-based landform segmentation of
LiDAR imagery. *International Journal of Applied Earth Observa-
tion and Geoinformation* 6(3), 261–270.
- McGarigal, K., Cushman, S. A., Neel, M. C., Ene, E., 2002.
FRAGSTATS: spatial pattern analysis program for categorical
maps. Tech. rep.
- Moller-Jensen, L., Kofie, R. Y., Yankson, P., 2005. Large-area urban
growth observationsa hierarchical kernel approach based on image
texture. *Geografisk Tidsskrift-Danish Journal of Geography* 105,
no. 2 (2005): 105(2), 39–47.
- Netzel, P., Stepinski, T. F., 2015. Pattern-based assessment of land
cover change on continental scale with application to NLCD 2001-
2006. *IEEE Transactions on Geoscience and Remote Sensing*
53(4), 1773–1781.
- Niesterowicz, J., Stepinski, T. F., Dec. 2013. Regionalization of multi-
categorical landscapes using machine vision methods. *Appl. Ge-
ogr.* 45, 250–258.
- Pontius, R. G., 2002. Statistical methods to partition effects of quan-
tity and location during comparison of categorical maps at multi-
ple resolutions. *Photogrammetric Engineering & Remote Sensing*
68(10), 10411049.
- R Core Team, 2013. *R: A Language and Environment for Statistical
Computing. R Foundation for Statistical Computing, Vienna, Aus-
tria.*
- Remmel, T. K., Csillag, F., 2006. Mutual information spectra for com-
paring categorical maps. *International Journal of Remote Sensing*
27(7), 14251452.
- Richards, J. A., 1999. *Remote sensing digital image analysis.*
Springer-Verlag, Berlin.
- Rubner, Y., Tomasi, C., Guibas, L. J., 2000. The earth mover’s dis-
tance as a metric for image retrieval. *International Journal of Com-
puter Vision* 40(2), 99–121.
- Shyu, C., Klaric, M., Scott, G., 2007. GeoIRIS: Geospatial infor-
mation retrieval and indexing system - Content mining, seman-
tics modeling, and complex queries. *IEEE Trans. Geosci. Remote
Sens.* 45 (4), 839–852.
- Steiniger, S., Hay, G. J., 2009. Free and open source geographic in-
formation tools for landscape ecology. *Ecological Informatics* 4(4),
183–195.
- Stepinski, T., Netzel, P., Jasiewicz, J., 2014. LandEx - A GeoWeb tool
for query and retrieval of spatial patterns in land cover datasets.
*IEEE Journal of Selected Topics in Applied Earth Observations
and Remote Sensing* 7(1), 257–266.
- Stepinski, T. F., Cohen, J. P., 2014. Comparing semantically-blind
and semantically-aware landscape similarity measures with appli-
cation to query-by-content and regionalization. *Ecological Infor-*

1120 matics 24, 69–77.
1121 Uuemaa, E., Antrop, M., Roosaare, J., Marja, R., Mander, b., 2009.
1122 Landscape metrics and indices: an overview of their use in land-
1123 scape research. *Living Rev. Landsc. Res.* 3 (1), 1–28.
1124 Vatsavai, R. R., 2013a. Gaussian multiple instance learning approach
1125 for mapping the slums of the world using very high resolution im-
1126 agery. In: In Proceedings of the 19th ACM SIGKDD international
1127 conference on Knowledge discovery and data mining. ACM, pp.
1128 1419–1426.
1129 Vatsavai, R. R., 2013b. Object based image classification: state of the
1130 art and computational challenges. In: In Proceedings of the 2nd
1131 ACM SIGSPATIAL International Workshop on Analytics for Big
1132 Geospatial Data. ACM, pp. 73–80.
1133 Williams, E. A., Wentz, E. A., Apr. 2008. Pattern Analysis Based on
1134 Type, Orientation, Size, and Shape. *Geogr. Anal.* 40 (2), 97–122.
1135 Zucker, S. W., 1976. Region growing: Childhood and adolescence.
1136 *Computer graphics and image processing* 5 (3), 382–399.