

CLASSIFYING PATTERNS OF LAND COVER USING MUTUAL INFORMATION AND CLUSTERING

Tomasz F. Stepinski

Lunar and Planetary Institute, 3600 Bay Area Blvd., Houston, TX 77058, USA

ABSTRACT

Collection of land use/land cover maps sharing common set of categories is classified on the basis of their patterns. Information-theoretic definition of similarity, based on the concept of mutual information is used to calculate multi-scale similarities between all map-pairs in the collection. Information contained in the derived similarity matrix is utilized for classification of all maps using a combination of visualization and agglomerative clustering techniques. The methodology is illustrated using a collection of 18 maps depicting land covers over major metropolitan areas in the United States. The cities are classified into five groups each having a distinct land cover pattern. In addition, all city-pairs in the collection are classified into five groups on the basis of similarity vector – a quantity that encapsulates map similarities at all available spatial scales.

Index Terms— Land cover, spatial patterns, pattern classification, urban environment

1. INTRODUCTION

Comparison between two maps of land use/land cover (LULC) is a fundamental task in remote sensing and geospatial data analysis with application to change detection, validation of models, and accuracy assessment [1, 2]. We propose to apply the methodology of quantifying agreement between pairs of LULC maps to a large collection of maps in order to classify landscapes into classes of LULC patterns. A pattern is a specific composition of LULC categories and their spatial arrangement in a given landscape; it represents a higher level abstraction of landscape than a single LULC category. For example, a LULC map of a city serves as means for visual assessment of spatial relations between its constituent categories, but it also defines a pattern - characteristic fingerprint of this particular city in terms of LULC. A collection of different cities can be grouped into classes on the basis of similarities between their patterns.

Traditional approaches to quantification of LULC patterns focus on either composition (histogram of LULC categories)

or spatial configuration (various pattern indicators). However, in the context of map comparison, it is sufficient to quantify only a *similarity measure* between two patterns without quantifying patterns themselves. We use a similarity measure [2] that simultaneously takes into account composition and configuration information. First, this similarity measure is applied to every pair of maps in a collection; the output of this step is a similarity matrix that encapsulates pairwise similarities between the maps. Second, we use a combination of clustering and visualization methods to translate the information contained in the similarity matrix into classification of landscapes (maps). The output of this step is a similarity map that depicts the overall structure of similarities between maps in the collection, and a classification of maps into characteristic LULC-pattern groups.

2. METHODS

We refer to a LULC map simply as a map and to LULC categories as colors. Distance between the two maps A and B , denoted as $d(A, B)$, is a symmetric function that increases with an increased dissimilarity between A and B . Similarity between maps A and B is given by $\text{sim}(A, B) = 1/(1 + d)$; the range of $\text{sim}(A, B)$ is between $\text{sim} = 0$, where the two maps are as distinct as possible, and $\text{sim} = 1$, where the two maps are identical.

Distance between the two maps is calculated using a method [2] which relies on the concept of normalized mutual information. The information-theoretic notion of mutual information [3] describes the amount of information (measured in units of bits) shared by stochastic processes. In application to map comparison the stochastic processes are natural and/or human actions that result in mapped distribution of colors. Mutual information is an absolute measure of information common to the two maps being compared; it can be transformed into a distance $d(A, B)$ through normalization. In [2] mutual information between a stochastic process of selecting a map (or a region in the map) and a stochastic process of selecting a color in the maps are used for calculating $d(A, B)$.

In order to expedite the calculations it is convenient (but not necessary) if both maps have square shapes with a size equal to $s = 2^p$ pixels. This allows to organize the information carried by the pixels into so-called quadtree (Q-tree)

This work was supported by the NSF under grant IIS-0812271 and by NASA under grant NNG06GE57G. The research was conducted at the Lunar and Planetary Institute, which is operated by the USRA under contract CAN-11NCC5-679 with NASA. This is LPI Contribution No.????.

data structure. The root of the tree is the entire map. The first four nodes contain data corresponding to the four quadrants (1 or NW, 2 or NE, 3 or SW, 4 or SE) of the map. Subsequent nodes carry data corresponding to regions of the map resulting from recursive subdivision of quadrants into four subquadrants. The terminal nodes (leaves) carry data pertaining to individual pixels. Single pixels are referred to by a list $Y = \{l_1, \dots, l_p\}$, where $l_i \in \{1, 2, 3, 4\}$, describing a path from the root of the Q-tree to an appropriate leaf. Larger regions (subquadrants) of the map are referred to by shorter paths.

A distribution of colors in two maps is given by a probability function $p(X, Y, Z)$, where X is a map variable ($X = A$ or $X = B$), Y is a spatial location variable (the Q-tree path), and Z is the color variable. Normalized mutual information between $p(X, Y)$ and $p(Z)$ is calculated as a proxy of “distance” between maps A and B . This quantity measures an average reduction in fraction of bits that are needed to convey (X, Y) if the distribution of Z is known – it may be interpreted as a relative information gain. Knowing that Z is distributed evenly among the two maps (the two maps have similar compositions) results in no information gain about (X, Y) , however, knowing that Z is distributed unevenly between the two maps (the two maps have dissimilar compositions) results in an information gain about (X, Y) . The Q-tree structure facilitates efficient calculations of mutual information between maps at multiple scales through simple adjustment of the length of Q-tree path Y ; calculations with length $(Y = 0)$ yield distance on the scale of the entire map (scale=1) whereas calculations with length $(Y = p)$ yield distance on the pixel scale (scale= $p+1$). The end result of mutual information calculations is a list $D = \{d_1, d_2, \dots, d_{p+1}\}$ of distances calculated at all the scales corresponding to all levels in the Q-tree structure; those distances may be converted to similarities using the formula given above. Repeating such calculation for all pairs of maps in a collection yields a similarity matrix; a single entry in the matrix is a list of similarities at various scales.

We use two complementary techniques to classify maps on the basis of similarity matrix. The primary technique is the agglomerative clustering. The advantage of using a clustering technique for classification of maps is its strict and precise result – the classes are accurately calculated as clusters. However, clustering does not offer an immediate insight into the overall structure of the map collection; the best number of clusters, their separations, and their within-cluster dispersions are not readily available from the results of a clustering algorithm. In order to address the shortcomings of clustering, we use data visualization as the complementary technique. The goal of visualization is to represent an overall structure of distances in the entire collection by a 2-dimensional graph (a similarity map). We use the Sammon’s map [4] technique for visualization because it requires only knowledge of distances between objects. The disadvantage of visualization

technique is its approximate character; the resultant planar graph reflects the structure of the collection only in an approximate fashion. Thus clustering and visualization techniques complement each other, using both improves the accuracy of classification.

3. DATA

In order to illustrate our method of classifying landscape patterns we extracted LULC maps for a collection of 18 large metropolitan areas in the United States from the National Land Cover Dataset 1992 (NLCD 1992) [5]. In order to expedite the calculations we degraded the original 30 m/pixel spatial resolution to 60 m/pixel and aggregated the original 21 LULC categories to only 10 (water, low intensity housing, high intensity housing, commercial/industrial/transportation, rock/sand/clay, forest, shrubland, grassland, agriculture/recreation, and wetland). We refer to the areas in the extracted maps as “cities.” Each map has a size 512×512 pixels and is centered at the center of a city. Thus, we classify LULC patterns of metropolitan areas contained in a square window with a side equal to ~ 30 km. Following cities are included: Houston TX (HOU), Chicago IL, (CHIC), Dallas TX (DAL), Denver CO (DEN), Los Angeles CA (LA), New York NY (NY), Philadelphia PA (PHIL), Seattle WA (SEA), Washington DC (WASH), San Francisco CA (SF), Phoenix AZ (PHX), Atlanta GA (ATL), Boston MA (BOS), Miami FL (MIA), Minneapolis MN (MSP), San Antonio TX (SAN), Saint Louis MO (STL), and Baltimore MD (BAL). We calculated mutual information-based distance and similarity measures for every city-pair in a collection at the scale of the entire map and at $\log_2 512 = 9$ smaller scales.

4. RESULTS

Fig. 1A shows an example of calculating similarity between LULC maps of two cities: HOU and NY. At the coarsest resolution (scale=1) the similarity is relatively small (the information gain is relatively large) due to significant differences in overall composition between the two maps. The similarity increases at spatial scales 2 to 4 indicating that, at these scales, the two cities have patterns that are more alike than the patterns for the entire maps. At finer spatial scales the similarity decreases implying significant differences at small-scale patterns. Similarity map of the entire collection can be calculated at any desired spatial scale by choosing only values of similarities at a specified scale level. The *overall* similarity between cities A and B , that takes into account all scales, can be calculated as, $\text{Sim}(A, B) = \sum_{i=1}^{p+1} V_i \text{sim}_i(A, B)$, where $\text{sim}_i(A, B)$ and V_i denote similarities and weights at scale i , respectively. Scale weights are subjective quantities that depends on the specific application.

Fig. 1B shows a similarity map of 18 cities constructed using a portion of similarity matrix corresponding to scale=1.

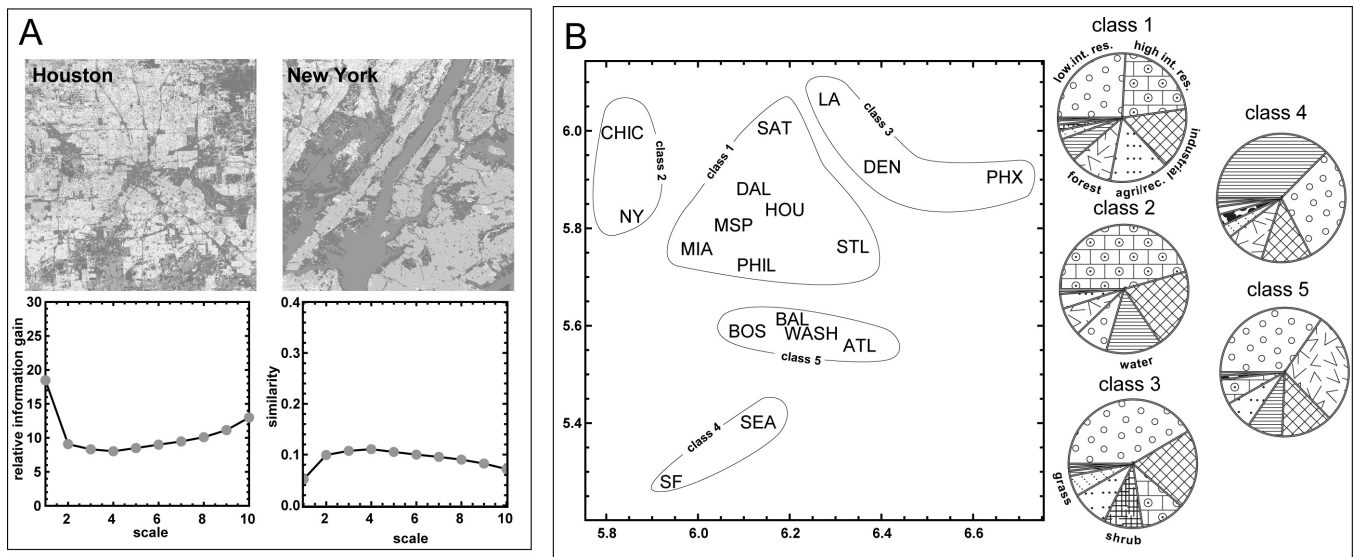


Fig. 1. (A) An example of similarity calculations between two LULC maps labeled as “Houston” and “New York.” Dependence on the spatial scale for both, relative information gain (distance) and similarity, are shown. (B) A similarity map of 18 cities; clustering-based classification into 5 classes is indicated by thin encirclements. Pie diagrams show characteristic composition of LULC categories in each class.)

Locations of cities on the similarity map are indicated by their abbreviated names; cities located close to each other on the map are similar whereas cities located far from each other on the map are dissimilar. The similarity map suggests that the collection of 18 cities may be divided into five clusters corresponding to distinct classes of LULC patterns. Agglomerative clustering provides memberships of these classes as indicated by encirclements on the similarity map. Pie diagrams located at the right side of Fig. 1 give characteristic compositions for the classes calculated as average compositions of their members. Class 1 (the largest class) groups cities dominated by residential areas (equal parts high and low density) with significant contributions from industrial and agricultural areas. Class 2, consisting of NY and CHIC, is dominated by high intensity residential areas with significant contribution from industrial areas and water. Class 3 groups cities dominated by low intensity residential areas with a significant contribution from industrial areas. Class 4, consisting of SF and SEA, is dominated by water and low intensity residential areas. Finally, class 5 is dominated by low intensity residential areas and forest. The membership in the classes or the classes themselves may change if the classification is based on similarities corresponding to a different scale level, or a composite similarity that weights contributions from similarities at different spatial scales (see above).

A pair of two cities is characterized by a “relationship” encapsulated by a list of similarities at all scales – a similarity vector. There are 153 pairs of cities in our collection, each pair is characterized by a unique similarity vector. Us-

ing the Euclidean distance as a measure of closeness between these vectors we can construct similarity map of relationships between city-pairs and cluster these relationships into groups congregating similar associations. Fig. 2A shows the similarity map of 153 city-pairs. Each city-pair is denoted by a symbol, pairs having similar relationships are located close to each other on the map and the pairs having different relationships are located far from each other on the map. We use agglomerative clustering to cluster all pairs into 5 groups corresponding to 5 different types of relationships. These groups are indicated by encirclements on the similarity map; in addition, pairs belonging to different groups are marked by different symbols. The membership of city-pairs in the groups is listed. Fig. 2B shows characteristic relationships (similarity vectors shown as curves) for all group calculated as average vectors from their members.

Group 1 corresponds to a “similarity” relationship. The high value of similarity at scale=1 implies similar overall composition of both cities. At finer spatial scales similarity decreases indicating differences in arrangement of LULC categories in the two cities. Group 5 (the largest group) corresponds to a “dissimilarity” relationship. The low value of similarity at scale=1 indicates significant differences in overall compositions of the cities. At finer scales the similarity remains small but increases slightly because somewhat more similar patterns are found within smaller regions of the two cities. Group 4 represents relationship similar to that given by group 5, but at somewhat larger levels of similarity; groups 4 and 5 could be combined into a single very large

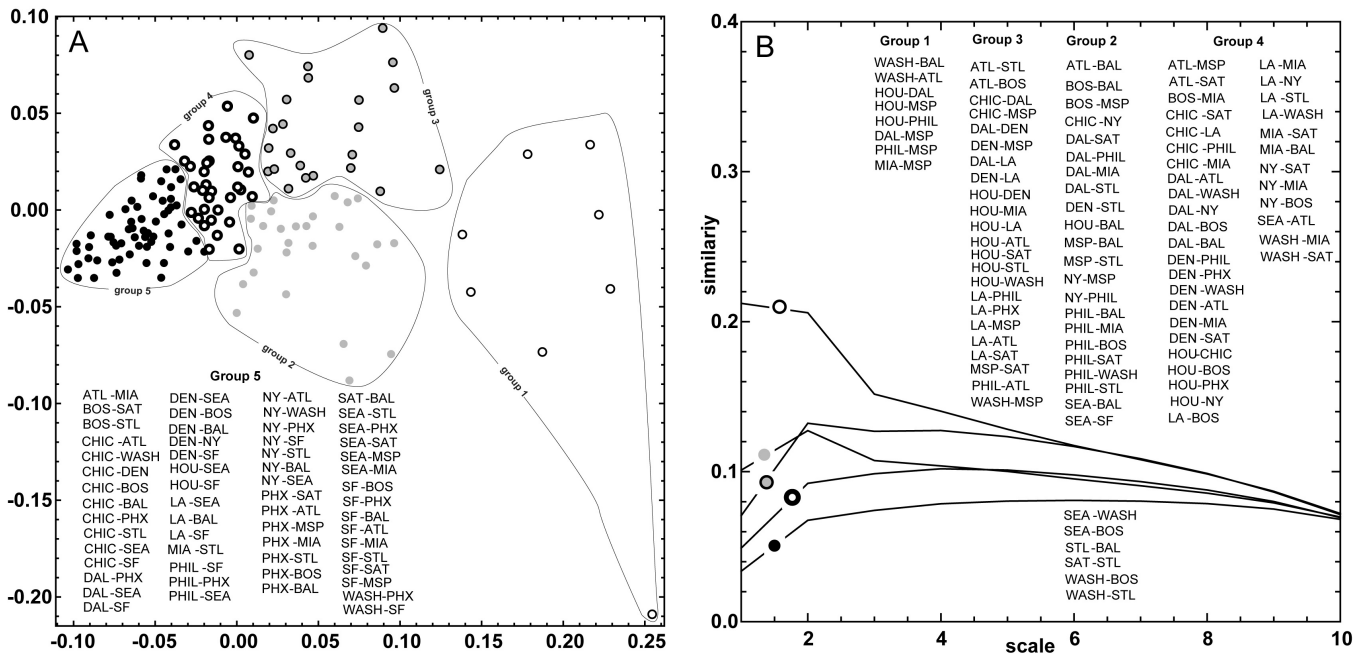


Fig. 2. (A) A similarity map of 153 pairs of cities; clustering-based classification into 5 groups is indicated and the membership in each group is listed. (B) Similarity curves characteristic for each group.

“dissimilarity” group. Group 3 corresponds to relationship that is defined by significantly different overall compositions but also by existence of matching patterns on smaller scales. Finally, group 2 corresponds to relationship defined by somewhat similar compositions at scales 1 and 2 and a gradual decrease of similarities at finer scales.

5. DISCUSSION

Presented method can find a broad application in all comparative studies of landscapes. One such application is to map landscape units larger than individual LULC categories – patterns of specific categories. For example, one of the cities in our collection may be divided into 64 tiles or sectors each having size of 64×64 pixels. The collection of 64 sectors can be classified into groups of similar patterns interpreted as downtown, suburbia, industry, etc. Assigning unique labels/colors to these groups results in a meta-map of a city bringing out urban units at higher level of generalization than the LULC categories. The method can also be applied to a collection of LULC maps derived at different times. This will help to catalog various trajectories of landscape change and to identify driving factors affecting landscape dynamics.

6. REFERENCES

- [1] Jr. R. G. Pontius, “Statistical methods to partition effects of quantity and location during comparison of categorical maps at multiple resolution,” *Photogram. Eng. Remote Sens.*, vol. 68(10), pp. 1041–1049, 2002.
- [2] T. K. Rimmel and F. Csillag, “Mutual information spectra for comparing categorical maps,” *Inter. Journal of Remote Sensing*, vol. 27(7), pp. 1425–1452, 2006.
- [3] C. E. Shannon, “A mathematical theory of communication,” *Bell System Technical Journal*, vol. 27, pp. 623–656, 1948.
- [4] Jr. J. W. Sammon, “A nonlinear mapping for data structure analysis,” *IEEE Transactions on computers*, vol. C-18(5), pp. 401409, 1968.
- [5] J. B. Vogelmann, S.M. Howard, L. Yang, C.R. Larson, B.K. Wylie, and N. Van-Driel, “Completion of the 1990s national land cover data set for the conterminous united states from landsat thematic mapper data and ancillary data sources,” *Photogram. Eng. Remote Sens.*, vol. 67, pp. 650–652, 2001.