

Quantifying the Visual Impact of Classification Boundaries in Choropleth Maps

Yifan Zhang and Ross Maciejewski

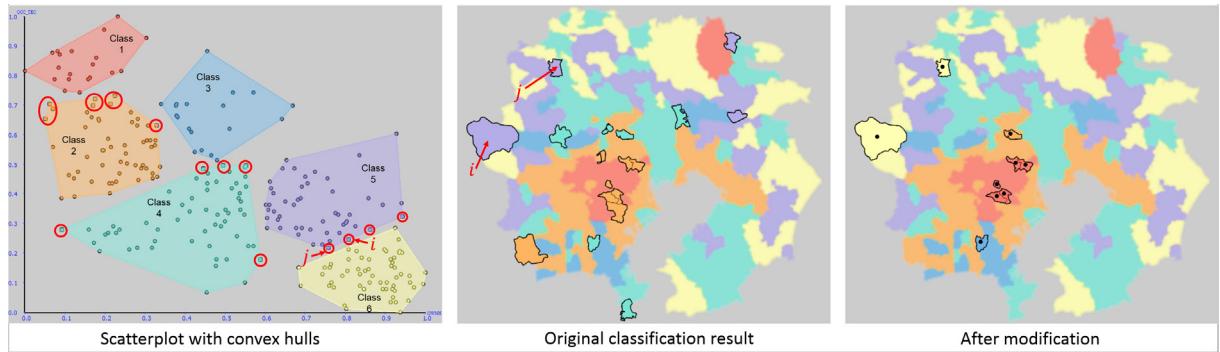


Fig. 1. An example of the visual effect of shifting classification boundaries. In this example, we create a choropleth map of housing factors in the Tokyo metro area (data captured in 1990 as part of a mortality study [40]). A two-dimensional k-means classification scheme was applied and elements near the classification boundaries are highlighted on the map and the scatterplot. The x-axis is the proportion of professional workers in an area and the y-axis is the proportion of the population that owns a house. The spatial units highlighted in the original classification result are all within some threshold value ($-0.03 \leq \tau \leq 0.0$) of the cluster boundary. If the cluster boundaries were to be modified, the result (after modification) would show different amounts of visual clustering of spatial units as demonstrated when comparing the two maps.

Abstract—One critical visual task when using choropleth maps is to identify spatial clusters in the data. If spatial units have the same color and are in the same neighborhood, this region can be visually identified as a spatial cluster. However, the choice of classification method used to create the choropleth map determines the visual output. The critical map elements in the classification scheme are those that lie near the classification boundary as those elements could potentially belong to different classes with a slight adjustment of the classification boundary. Thus, these elements have the most potential to impact the visual features (i.e., spatial clusters) that occur in the choropleth map. We present a methodology to enable analysts and designers to identify spatial regions where the visual appearance may be the result of spurious data artifacts. The proposed methodology automatically detects the critical boundary cases that can impact the overall visual presentation of the choropleth map using a classification metric of cluster stability. The map elements that belong to a critical boundary case are then automatically assessed to quantify the visual impact of classification edge effects. Our results demonstrate the impact of boundary elements on the resulting visualization and suggest that special attention should be given to these elements during map design.

Index Terms—Choropleth, Classification, Visualization, Geodemographics, Geovisualization

1 INTRODUCTION

One of the most common methods of visualizing spatially referenced data is the choropleth map. Choropleth maps are based on data aggregated over defined areal units (country, ZIP code, etc.), and one of the first design choices in creating a choropleth map is determining which range of data values should be associated with which color. This step is typically referred to as classification, and a variety of class interval selection/binning methods (e.g., quantile, equal interval, standard deviation, natural breaks [30], minimum boundary error [15], and genetic binning [5]) have been developed. The selection of the class interval has a major impact on the visual appearance of the map [18]. Ideally, regions that are alike under a given statistical measure will appear as the same color on a map; however, due to the nature of interval selection, there can be map elements where the statistical measure falls near the

classification boundary. Depending on the spatial positioning in the map, switching a map element from one class to another could result in large changes in the appearance of visual clustering in the map. For example, in Fig. 1, several elements are on the boundary between two classes (Yellow and Purple) as well as others. Here, we show that if the classification boundary was slightly shifted, the area under analysis would visually appear to have a larger amount of spatial clustering than if the boundary were to remain as chosen. Given unit i and j 's spatial proximity to other elements in the Yellow class, shifting the Purple-Yellow classification boundary slightly to incorporate these units into the Yellow class may not be unreasonable.

This issue is further compounded in more complicated maps where the classification is not being done for one or two variables but rather as a combination of multiple variables. Such multivariate classification is common practice in many areas, such as geodemographic profiling [6, 46], ecological area selection [25], and epidemiology [7], and involves creating intervals based on multiple statistical measures. This multivariate classification typically employs various data mining and machine learning classification methods (e.g., k-means [36, 43], self-organizing maps [23], hierarchical clustering [13, 22]), and the class intervals in the choropleth map no longer belong to a single data range but rather they belong to a more complex combination of relationships

• Yifan Zhang and Ross Maciejewski are with Arizona State University. For questions, e-mail: rmacieje@asu.edu.

Manuscript received 31 Mar. 2016; accepted 1 Aug. 2016. Date of publication 15 Aug. 2016; date of current version 23 Oct. 2016.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below.
Digital Object Identifier no. 10.1109/TVCG.2016.2598541

between variables. Thus, the issue of boundary cases (i.e., map elements that could potentially belong to multiple classes) is critical due to the fact that small changes to the classification of a boundary element can result in a large change in the visual appearance of the choropleth map.

Here, it is critical to note that the concept of map classification is a matter of reducing precision in order to provide a compact visual representation. What is being done is to group similar data elements into classes and project these on the map. Spatial analysis will be done on the actual data values in order to explore relationships in the data. However, the visual display of the data still depends on this reduced precision. As such, a map classification can imply to a viewer that spatial relationships exist (whether intentionally or unintentionally). Thus, we need to fully understand the visual implications that such classifiers have. In this paper, we utilize a cluster measure to identify map elements that could belong to multiple class intervals as a result of multivariate classification as done in prior studies [8, 50, 52]. Once these elements are identified, we develop a novel metric based on a global measurement of spatial clustering (Moran's I [39]) to quantify the visual change that will occur if a boundary element is relabeled. Based on this metric, we identify map elements that warrant a closer inspection as a reclassification of the elements will result in more (or less) spatial association being visually present in the map. We believe that the identification of such map elements will better inform map design by enabling analysts or designers to identify spatial regions where the visual appearance may be the result of spurious data artifacts, or to identify spatial regions that should potentially be reclassified due to their regional spatial association. Demos, datasets and code related to the proposed methodology can be found at: <https://vader.dtn.asu.edu:8443/BoundaryElements/>.

2 RELATED WORK

Our work focuses on quantifying the effect of the changes caused by modifying the class labels of boundary elements. In this section we review the work in map classification, spatial association, map comparison and direct manipulation related to classification boundaries.

2.1 Map Classification

The goal of a classification scheme is to group similar observations and split dissimilar observations to simplify and clarify the message of the map. For univariate data, the simplest methods include quantile, equal interval, and standard deviation [35]. More complex methods have been proposed since the early 1960s. For example, Jenks developed natural breaks, which seeks to reduce the variance within classes and maximize the variance between classes [30]. Scripter presented nested means [45] that calculates intervals for statistical maps by repeatedly deriving and using the arithmetic mean to divide a numerical array. Cromley [15] proposed a minimum boundary error method that maximizes spatial similarity among contiguous units in the same class interval, and Armstrong [5] developed a genetic binning scheme that creates optimal classifications with respect to multiple criteria (e.g., number-line relationships, fragmentation).

The most important part of map classification is how to choose the breaks or class boundaries. Evans [18] categorized sixteen class-interval systems and suggested that class intervals should be selected according to the overall shape of the data distribution. Brewer et al. [10] compared seven map classification methods with fifty-six subjects in a two-part experiment to determine which classifications are most suitable for epidemiological rate maps. Sun et al. [48] proposed a heuristic classification approach that utilizes the class separability concept and other classification criteria. They compared their approach to other classification methods based on element separability; however, visual changes in the map appearance that could occur due to slight shifts in classification boundaries have not, to our knowledge, been fully addressed. In fact, most of the previously mentioned studies focus on single variable classification methods where the statistical distributions can be easily plotted and explored.

Multivariate map classification, on the other hand, typically involves classification over several variables using various data mining and ma-

chine learning classification methods (e.g., k-means and self-organizing maps) that are often used as black box methods. By assigning a color to each label/ class/ category in the clustering result, a choropleth map is generated. Here, the multivariate statistical distributions become visually complex, and designers may simply default to the base parameters. Such multivariate classification methods are heavily utilized in demographics classification. For example, Vickers and Rees created the United Kingdom National Statistics Output Area Classification (OAC) [49], which is an open geodemographic classification with a hierarchical structure of 7 super-groups, 21 groups and 52 subgroups. Because of the efficiency and simplicity [28], k-means clustering remains the core algorithm for the computation of geodemographic classifications [27]. Therefore we use k-means as our default multivariate classification method; however, our findings can easily be extended to other classification methods.

2.2 Spatial Association

One of the main uses of choropleth maps is enabling analysts to mentally assess spatial associations. In this paper, we focus on how a modification of elements near the classification boundaries could potentially change the observed spatial association. To quantify this, we focus on statistical measures of spatial association, particularly spatial autocorrelation, which is a statistical measure of how spatial units are associated. A variety of methods for spatial autocorrelation have been developed over the past decades. The most popular metrics for global spatial autocorrelation include the Join count statistic, which was developed mainly for binary variables based on the probability of a unit having neighboring units of the same class [14], Moran's I [39], which considers pairwise products of deviations, Geary's C [19], which considers pairwise squared differences, and Getis-Ord General G [20], which is mostly used for hotspot detection. Global Moran's I, Geary's C and Getis-Ord G are developed for measuring autocorrelation in continuous variables, and the main difference between these methods is the measures of value similarity used. These methods can be generalized into cross-product statistics [29], and local versions of these metrics have been further developed (e.g., LISA [4]).

While methods for measuring autocorrelation in continuous variables are critical, our focus is on the autocorrelation between class labels, which are categorical variables. In order to compute categorical spatial association, Join count statistics have been used and are often applied to the k-color cases [16]. For example, Boots [9] developed a procedure for extending local statistics to categorical spatial data based on the composition and configuration characteristics of categorical data. Our work leverages these spatial autocorrelation metrics as a means of identifying critical points with regards to the change of visual appearance. In this way, we can quantify the visual change that will occur if classification boundaries are changed.

2.3 Map Comparison

The major metric of assessing changes based on re-classification of data in a choropleth map is through visual inspection and comparison. Research in comparison of choropleth maps has focused on using statistics (e.g., intercorrelation [33], relative blackness [34]) to visually compare spatial distributions. Olson [41] examines the effects of class interval systems via the visual judgment of the correlation between pairs of choropleth maps. Lloyd et al. [34] shows that decisions on the similarity of maps appears to be influenced by both the similarity of the spatial distributions and the relative blackness of the maps. Olson [42] also explores the issue of map pattern complexity regarding several statistics including rank autocorrelation, average differences, and weighted proportions. Xiao and Armstrong [51] developed an evolutionary algorithm that allows users to explore spatial patterns in terms of their visual correlation. Our work differs from previous research in that we focus on the effect of spatial association changes due to boundary modifications rather than conventional visual correlation; however, these measures that link similarity to appearance are directly applicable and could be extended in future work as yet another quantification method for assessing the impact of boundary cases in map design. Based on previous literature, it is clear that the size of a region

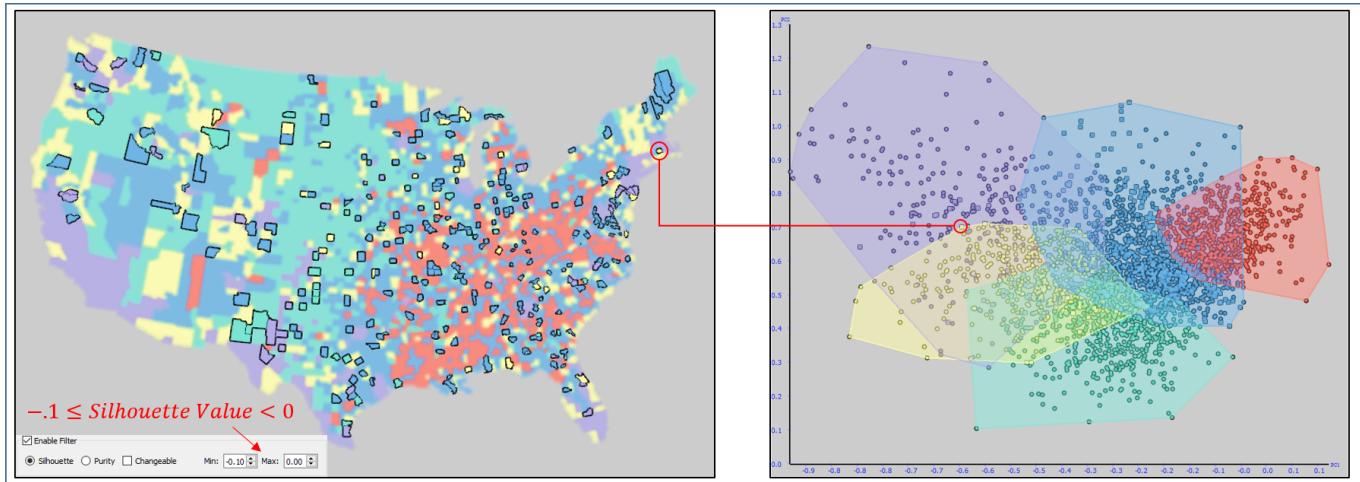


Fig. 2. An example of k-means clustering ($k = 5$) using the US Census data variables “Education above bachelor’s degree,” “Mean time to travel to work,” and “Foreign born person.” Boundary elements with a silhouette coefficient from $-.1 \leq \tau \leq 0$ are highlighted. A projection in the geographical space is shown on the left, and a principle component projection of the first two principle components is on the right. We annotate one boundary element with a red circle and show the relationship between the two views.

(number of units in a region) directly impacts the implied spatial associations [31] of a map. Thus, as a contiguous region adds more elements on the map, the size of the region may also imply more importance, and work by Haroz and Whitney [26] illustrated that grouped arrangements of elements directly influence visual search and subitizing tasks. As such, our metric directly ties to the contiguity of units as well as the change in the size of a region; however, other metrics such as shape and overall proximity should be considered for future exploration.

2.4 Direct Manipulation & Classification Boundaries

Given that multivariate map classification utilizes a variety of data mining and machine learning algorithms, it is critical to note that many methods have been developed to help analysts interactively explore and manipulate cluster structure. The goal is to utilize domain knowledge to refine classifications and reduce classification errors. Such work is directly applicable to classification methods for generating choropleth maps. Work in this area includes the VISTA system [11], which was developed to help domain experts validate and refine cluster structures through interactive feedback. VISTA allows users to mark the visual boundaries between clusters and refine the algorithmic result if applicable. Chen and Liu [12] developed iVIBRATE as an interactive machine learning tool, which allows users to iteratively modify the clustering process. iVIBRATE consists of a visual cluster rendering component and an adaptive labeling subsystem, and Andrienko et al. [2] developed methods for analyst-guided clustering of large collections of trajectories by combining clustering and classification together through an interactive interface. In this paper, we focus on the identification of class elements near boundaries and enable direct manipulation for relabeling class elements. While our focus is more on the resulting changes in the visual output, the identification of elements that impact the visual output can be used as measures of importance to direct analysts’ attention to elements that require further inspection.

3 THE VISUAL IMPACT OF BOUNDARY ELEMENTS

The analysis and understanding of spatial patterns is essential to all subfields of geography, and the visual representation of spatial patterns is greatly affected by the choice of classification boundaries. Applying an inappropriate classification may create false patterns or lead to misinterpretation of the resulting map and choosing an appropriate data classification scheme for map generation can be difficult [32]. In fact, Klipper et al. [31] found that subjects (both expert and non-expert) seemed to base their notion of spatial significance of a map on the number of cells of a particular color. Thus, the choice of class for boundary elements can directly impact the perceived spatial significance as these classes directly lend themselves to the cell count. Given these known

issues, it is critical to evaluate the map classification design prior to presentation. By identifying elements that lie near a classification boundary, we can quantify the visual impact that shifting the boundary will have on the map. In this way, users can explore and modify their classification design scheme or highlight problematic elements that may be contributing to spurious visual effects.

Previous work has also focused on exploring elements near classification boundaries, for example, Egbert and Slocum created ExploreMap [17], which showed zones that are near classification boundaries in the univariate case. Andrienko et al. [3] developed the Descartes [1] system to provide additional statistics (min, max, etc.) about available attributes for each class in order to help analysts understand the classification relationship, and recent work by Slingsby et al. [47] shows class and distance to class centroid for the results of a multivariate classification for all locations concurrently. While distance to the centroid can provide insight into the compactness of a cluster, it is not necessarily a measure of the stability of the classification of a map element. Our work extends beyond previous work and explores the use of other indicators of cluster stability (as opposed to distance to the centroid) and defines the visual impact of class boundaries. By identifying elements that are potentially unstable in a cluster, we are then able to assess how a change of the element’s classification will impact the resultant visualization. As such, a discrete method for identifying only elements that impact the visual output provides designers with new information not found in previous methods.

3.1 Identification of Boundary Elements

A variety of metrics exist for characterizing the stability of the results of multivariate clustering. One such metric is the silhouette coefficient [44], which is used to define the separation distance of elements between the resulting clusters. The silhouette coefficient is defined as:

$$S(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}},$$

where $a(i)$ is the average dissimilarity of object i with all other objects within the same cluster, and $b(i)$ is the lowest average dissimilarity of object i to any other cluster in which object i is not a member. $S(i)$ is bounded by $-1 \leq S(i) \leq 1$. $S(i) = 1$ indicates that element i is very far away from all other clusters and so is most likely classified correctly. $S(i) = 0$ indicates that element i is near (or on) the decision boundary of a cluster (meaning it could potentially be reclassified), and $S(i) < 0$ indicates that element i is likely misclassified. We leverage this coefficient as a means of assessing the boundary elements between clusters. A range, τ for the silhouette coefficient is interactively chosen,

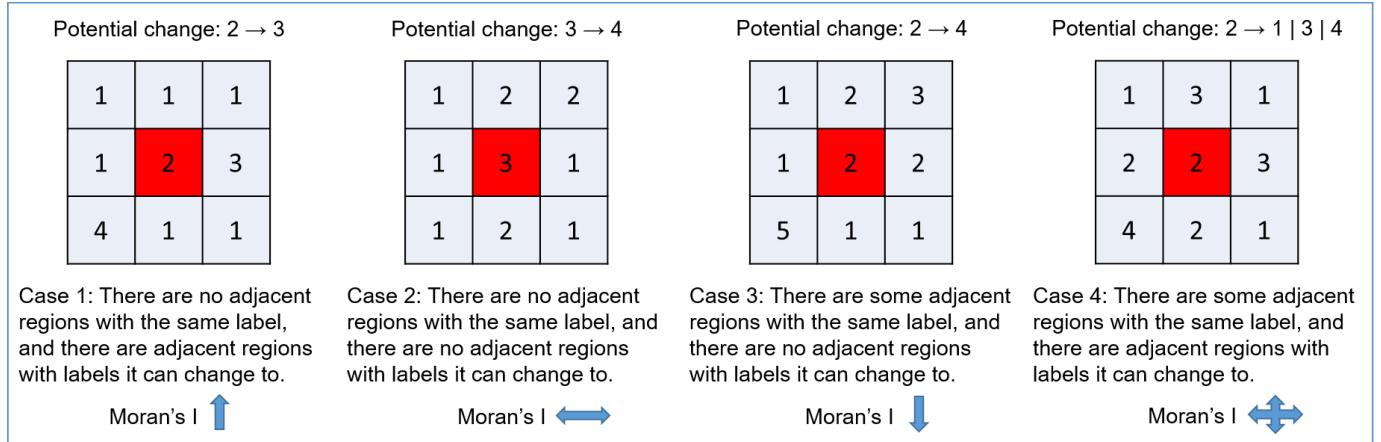


Fig. 3. Four spatial cases and the effects of changing a single unit.

and map elements with values satisfying τ are highlighted on the map. For example, Fig. 2 shows a demographic clustering of three US Census variables (“Education above bachelor’s degree”, “Mean time travel to work”, and “Foreign born person”) with $k = 5$.

In order to identify which boundary an element i is associated with during the computation of the silhouette value, the proximity of i to each cluster is stored. Then, an ordered list of classes are assigned to i based on the chosen silhouette range, τ . The ordered list represents all the potential classes that i could be reclassified as.

In addition to the numerical metrics for identifying the points near boundaries, we also provide a principle component analysis (PCA) scatterplot view (Fig. 2 right) to enable users to visually inspect boundary conditions of projected clusters. By projecting the k-means clusters into a 2D space, users can have a generalized overview of how the clusters are distributed. As the silhouette range τ is changed, boundary elements are indicated by changing shape from circles to squares in the PCA scatterplot.

3.2 Indicators of Spatial Association

While the identification of boundary elements is critical, the silhouette coefficient does not provide information on whether changing these elements will impact that visual spatial associations on the map. In order to identify the visual impact of shifting an element class, we first define the types of spatial association that can be observed:

- *Clustered*: Map elements with the same class are contiguous in geographic space, as indicated by positive measures of spatial autocorrelation in Moran’s I.
- *Dispersed*: Map elements with different classes (but with a repeated pattern) are contiguous in geographic space, as indicated by negative spatial autocorrelation in Moran’s I, an example of such a pattern would be a checkerboard.
- *Random*: Map element classes are randomly distributed on the map, as is indicated by a Moran’s I near zero, i.e., the distribution of regions with similar properties is unspecified/random in geographic space.

Each type of pattern is associated with a description of the visual appearance of the map, and these spatial association patterns are typically defined and tested using spatial autocorrelation. Spatial autocorrelation is often used with p-value, z-score, and resampling methods to indicate the significance level of the tendency of spatial clustering in a map. Our goal is to adapt an indicator of spatial association to quantify the visual change that may occur in a choropleth map as an element’s class is altered.

Many indicators for spatial association exist (e.g., join count statistics [14], Geary’s C [19], Moran’s I [39], Getis-Ord General G [20]).

However, these statistics are all special cases of cross-product statistics [14, 29]. Moran’s I [39] is perhaps the most well-known and widely used measure of spatial autocorrelation. Moran’s I is defined as:

$$I = \frac{N}{\sum_i \sum_j w_{ij}} \frac{\sum_i \sum_j w_{ij} (x_i - \bar{X})(x_j - \bar{X})}{\sum_i (x_i - \bar{X})^2}, \quad (1)$$

where N is the number of spatial units indexed by i and j , x is the variable of interest, \bar{X} is the mean of x , and w_{ij} is an element of a matrix of spatial weights.

Unfortunately, Moran’s I is designed for continuous variables. Since the visual appearance of the map relates solely to the final class, we need a metric that can be applied to categorical data values. As such, we modify the Moran’s I measure to provide a metric of spatial autocorrelation based on the class. To do this, we need to redefine the variables in Equation (1). x_i is now defined as a vector (c_1, c_2, \dots, c_n) , where n is the number of clusters and c_n is a binary value, 0 or 1, such that if element i belongs to cluster 1, then $c_1 = 1$ otherwise, $c_1 = 0$. Then \bar{X} will be the average of all vectors x_i , and a modified global Moran’s I can be calculated to evaluate the spatial association of classes. In this paper, we utilize the Queen contiguity for defining the spatial weights matrix. And $w_{i,j} = 1$ for all Queen contiguous neighbors in our implementation. While the choice of the spatial weights matrix will impact the calculation, the application is generalizable to any spatial weights choice. Note that the change of the definition of X_i is due to the fact that we are applying Moran’s I over the results of the map classification (instead of using the statistical measure of a county, we are using the class of the county). Since the county value is categorical, we cannot use the continuous formulation for X_i . While Join counts are appropriate for categorical data, our choice of using Moran’s I is due to the fact that it is one of the most commonly used measures of spatial autocorrelation. Furthermore, an application of join count would require treating the K-color case as a binary case (1 color and K-1 colors). This would result in multiple computations of the join count and added complexity to the proposed methodology.

3.3 Categorizing the Effects of Reclassifying

Once a measure for the spatial association of the classes is defined, the next step is to determine the cases in which altering a class will impact the visual spatial association. We identify four potential spatial arrangements for elements on a choropleth map, Fig. 3. Based on these arrangements, we then define the value change in our modified Moran’s I that would result in a change of the classification.

Case 1: The spatial unit under analysis, i , is spatially contiguous only to units with different classes. The position of i in the classification space is such that it lies near the class boundary of one or more spatially contiguous units. In this case, if i was reclassified, the spatial association will increase. This is illustrated in Fig. 3

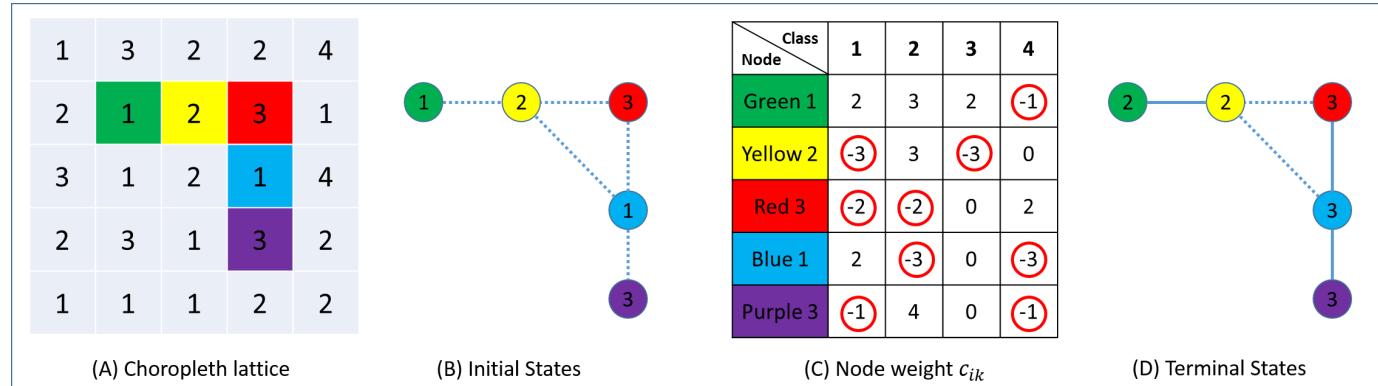


Fig. 4. An example of adjacent changeable regions. Dashed lines represent the contiguity and solid lines represent the co-effect. Non-negative node weight indicates that class k is reachable by i . From the initial states to the terminal states, three co-effect connections have been established.

(Case 1). Here, unit i is the red square and belongs to class 2. This element lies near the boundary of class 2 and class 3. If the class of i were to change from 2 to 3, an increase in visual clustering could be observed and the spatial association value would increase.

Case 2: The spatial unit under analysis, i , is spatially contiguous only to units with different classes. The position of i in the classification space is such that it does not lie near the class boundary of any spatially contiguous units. In this case, if i was reclassified, there would be no change in the spatial association. This is illustrated in Fig. 3 (Case 2). Here, unit i is the red square and belongs to class 3. None of its neighbors share the same class, thus i does not add to any visual cluster. i lies on the boundary of class 3 and class 4; however, changing i 's class to 4 does not result in i visually combining with other spatially contiguous regions, thus there is no change in the spatial association metric.

Case 3: The spatial unit under analysis, i , is spatially contiguous to some (or all) units that share the same class. The position of i in the classification space is such that it does not lie near the class boundary of any other spatially contiguous units. In this case, if i was reclassified, the spatial association will decrease. This is illustrated in Fig. 3 (Case 3). Here, unit i is the red square and belongs to class 2. Several of its neighbors share the same class, thus forming a small region that will visually appear to be clustered. While i does lie near the boundary of class 2 and class 4, there are no spatially contiguous elements belonging to class 4. As such, if i were to be reclassified, the size of the region containing elements belonging to class 2 would decrease, and no other region in this scenario would add i to their spatial grouping. As such, the visual clustering would decrease, resulting in a lower spatial association value.

Case 4: The spatial unit under analysis, i , has a class that lies near a classification boundary and is spatially contiguous to some units that share the same class. The position of i in the classification space is such that it does lie near the class boundary of other spatially contiguous units. In this case, if i were to be reclassified, the change in spatial association could be positive, negative, or neutral dependent on the number of contiguous units (and their contiguous units) that have the same class as i . This is illustrated in Fig. 3 (Case 4). Here, unit i is the red square and belongs to class 2. Several of its neighbors share the same class, thus forming a small region that will visually appear to be clustered. However, i lies on the boundary of class 2 and class 4 and is spatially contiguous to other regions belonging to class 4. If i were to be reclassified, the size of the region containing elements belonging to class 2 would decrease; however, the size of the region containing elements belong to class 4 would increase. As

such, the modified Moran's I would need to be recalculated for the entire map to determine the net change in spatial association.

While Cases 1-3 are straightforward to identify, Case 4 is perhaps the more common case in choropleth map design. Thus, for a unit i in Case 4, we define the number of regions that belong to the same cluster as i in its surrounding area as p_i . The number of regions that belong to a cluster that i can change to in its surrounding area as q_i . The effect on the spatial association after i is changed is based on the number of surrounding units that i can change to and is proportional to $q_i - p_i$. Fig. 3 only considers the effect of a single changeable unit, we extend this to more complex situations (Fig. 4(A)) in which several contiguous regions could change, resulting in a cascade of visual clustering patterns.

Theorem 1. *If the potential changeable regions are not adjacent, then their effects on the spatial association are separate/independent.*

By inspection, one can observe that if spatial units that are identified as being near class boundaries are non-adjacent, then the effect of modifying their classes will be independent. This can be observed in Equation (1) where units that are not adjacent will have an entry in the spatial weights matrix $w_{ij} = 0$ making the resulting calculations independent from one another.

Once independence is established, we can identify all spatial units that fall into Cases 1-4. Then, we can consider the situation where several changeable regions are adjacent, meaning that a change of class in one region will affect the visual clustering (i.e., the value of p and q) of another changeable region. In this case, we have:

Theorem 2. *The effect of the change (EOC) only depends on the initial states and the terminal states of the changeable regions.*

Thus, the measurement of spatial association remains the same as long as the final states of those changeable regions stay the same. We generalize the effect of the changes as:

$$EOC_{\xi} = \underbrace{\sum_{i \in \xi} (q'_i - p'_i)}_{A} + \underbrace{\frac{1}{2} \sum_{i \in \xi} \sum_{j \in \xi, j \neq i} w_{ij} (i_t \& j_t - i_s \& j_s)}_{B}, \quad (2a)$$

$$i_t \& j_t = \begin{cases} 1 & \text{if } i_t = j_t \\ 0 & \text{if } i_t \neq j_t \end{cases} \quad (2b)$$

where ξ is the set of changeable units, q'_i, p'_i are similar to q_i, p_i but exclude the other changeable units. w_{ij} is the spatial weight between spatial units i and j . i_t and j_t are the terminal states (classes) of regions i and j respectively, and i_s and j_s are the starting classes of regions i and j respectively. Here the effect of the changes can be broken into the total separated effect caused by all of the changeable regions (Equation (2a) A) and the total co-effect among those changeable regions

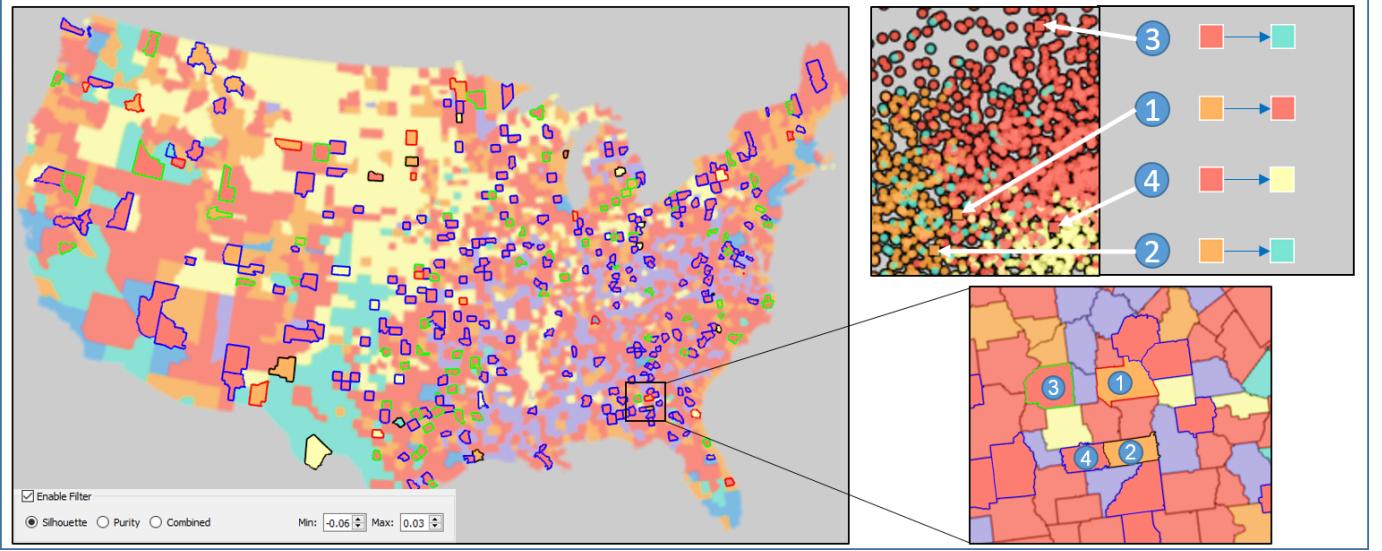


Fig. 5. A k-means classification of US census variables illustrates boundary elements and their corresponding cases from Fig. 3. Here the Red outline represents Case 1, the Black outline represents Case 2, the Green outline represents Case 3, and the Blue outline represents Case 4.

(Equation (2a) B). Note that the co-effect is divided by 2 because i and j are symmetric and would double the effect.

We can maximize EOC in Equation (2a) to determine the set of classes that will create the largest visual clustering in the map. This problem can be solved by maximizing the modified Moran's I in the terminal class state of a unit. We note that this may not be a desirable effect as this could introduce spurious patterns into the map; however, elements near classification boundaries need to be inspected and the EOC can be used as a metric for defining which elements could have the largest potential change on the visual output (which is the overall goal of this work). First, it is assumed that there exists a group of contiguous spatial units that lie near classification boundaries (Fig. 4(A)). Each unit can be altered to a certain class with a known weight. The weight is set to be the number of neighboring units that share the same class. In practice, for the class that a unit i cannot change to, the weight is set to $-\sum_{j \in \xi, j \neq i} w_{i,j}$ (see the red circle in Fig. 4(C)). By setting the weight to $-\sum_{j \in \xi, j \neq i} w_{i,j}$, we neutralize the possible co-effects and guarantee that a unit cannot change into an unreachable class. If the two adjacent units have the same class, an edge will be established with a given weight. For simplicity, the weight of the edge is unified to 1 when the spatial weight $w_{i,j}$ is 1. Finally, this can be formulated as a maximization problem where the nodes need to be classified such that the overall weight of the nodes and edges is maximized. This can be further defined as an integer linear programming (ILP) problem. Given a graph $G = (V, E)$ with n nodes and each node has m choices of classes, we introduce binary variables x_{ik} ($i = 1, \dots, n$, and $k = 1, \dots, m$) to indicate whether node i has been classified as class k . The weights $c_{ik} \in \mathbb{R}$ are given for each x_{ik} , and variables $y_e, e \in E$ indicate whether edge e is valid based on if its two nodes have been classified in the same class (Fig. 4(D)). The resulting ILP can be formulated as:

$$\max \quad \sum_{i=1}^n \sum_{k=1}^m c_{ik} x_{ik} + \sum_{e \in E} y_e \quad (3a)$$

$$\text{s.t.} \quad \sum_{k=1}^m x_{ik} = 1 \quad i = 1, \dots, n \quad (3b)$$

$$2y_e - x_{ik} - x_{jk} \leq 0 \quad e = (i, j) \in E, k = 1, \dots, m \quad (3c)$$

$$x_{ik} + x_{jk} - y_e \leq 1 \quad e = (i, j) \in E, k = 1, \dots, m \quad (3d)$$

$$0 \leq x_{ik}, y_e \leq 1 \quad (3e)$$

$$x_{ik}, y_e \in \mathbb{Z}. \quad (3f)$$

Here Equation (3c) constrains two nodes of a valid edge to be in the same class and Equation (3d) constrains an invalid edge to not have

two nodes in the same class. By solving this ILP we can identify the terminal states that maximize the Moran's I, the same formulation can also be used to minimize the Moran's I. The problem of finding the maximum possible value is similar to the Maximum Edge-Weighted Clique Problem (MEWCP) [37], which is a known NP-Hard problem. Therefore the problem of calculating the maximum EOC is also an NP-Hard problem (i.e., there is no general solution that can find the optimized value in polynomial time). Efficient algorithms for the MEWCP, such as heuristics approximation, may be modified and applied to solve this ILP. However, in practice, we find the number of adjacent nodes and the number of class choices are relatively small (traditional choropleth map design rules of thumb limit the number of classes to be less than 9). Thus, we implement a brute force solution to compute all possible values in our framework. During this computation, our framework stores the configuration of the classes that would maximize or minimize the current spatial association. Fig. 5 shows a multi-dimensional classification of demographic data in the United States and we highlight county boundaries based on their correspondence to the cases of Fig. 3 for illustrative purpose.

Note that for a connected component with n spatial regions and m possible class choices, the time complexity is $O(m^n)$. As such, our computations are heavily dependent on the number of boundary elements that have been identified. In practice, the largest dataset we have tested is the county map for the continental United States, which contains over 3100 spatial regions. A 6-class map classification with Silhouette value between -0.4 to 0.1 identifies approximately 600 regions that lie on class boundaries. Among these units, there are about 50 connected components and the largest connected component contains 20 regions. Our experiments used an Intel Core i7-3630QM Quadcore 2.40 Gz, and we found that the Min or Max EOC calculation time was less than 1 second.

3.4 Unit Size as a Function of Visual Change

While our proposed metric for quantifying the impact of visual change takes into account classes, perceptual studies have also shown that the size of the map units is a primary driver behind the patterns that users observe. As noted by Haklay [24], a thematic map created using spatial units that vary in shape and size leads the user into thinking that the larger areas are more significant because they have a bigger visual impact than the smaller areas. Seonggook [38] proposed the concept of gross change detection and verified that different spatial distributions between two adjacent choropleth maps may lead users to under- or over-estimate the gross change in the map, which implied that the spatial distribution of change should be considered. As such,

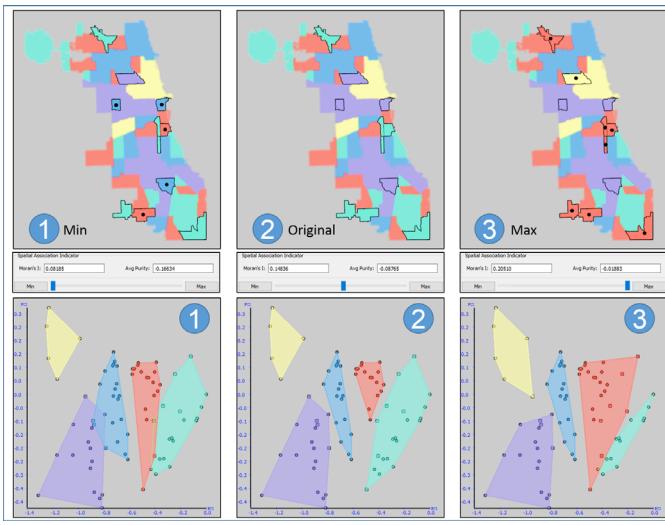


Fig. 6. Minimizing and maximizing the EOC of elements near the classification boundary using criminal incident reports in Chicago, IL. 1 - Minimizing the EOC, thus creating more visual heterogeneity, highlighted elements are those near the classification boundary. 2 - The initial k-means classification with only the changeable elements highlighted. 3 - Maximizing the EOC, thus creating more visual spatial clustering.

the size of the region should be considered when quantifying the visual impact of the classification changes. Goldsberry and Battersby [21] introduced the magnitude of change (MOC) to quantify the graphical change between choropleth map pairs for animated choropleth maps. MOC is applicable to both object-oriented and pixel-based measures, and we extend our EOC measure to consider the size of the map element with a final metric for quantifying the impact of boundary effects on the visual spatial association in choropleth maps. The metric is a simple multiplication to derive the visual impact of changes (VIOC) and is defined as:

$$VIOC_{\xi} = \sum_{i \in \xi} \left(\frac{s_i EOC_i}{S} \right), \quad (4)$$

where s_i is the area size of the i th region (in pixels) and S is the overall area size of the map (in pixels). This accounts for the proportional physical change of the choropleth map under different resolutions.

3.5 Summarizing the Visual Impact

Once these metrics are defined, we can now identify units on the map that could potentially be modified to change the visual appearance of spatial association. While there are methods for specifically identifying statistically significant spatial associations on a map, the majority of choropleth maps are presented with no underlying analysis of spatial association. Instead, they are presented in the wild and left solely for visual interpretation. By being able to quantify potentially spurious elements on a map, new designs could be considered where the elements could be blurred, highlighted or reclassified to another separate class in order to try and insure that patterns being seen are what was intended by the map designer (of course we recognize that the intent of the designer could have been to mislead). Thus, our method could be summarized into the following steps:

1. Choose a classification method and classify the dataset of interest
2. Calculate the silhouette value for all elements in the dataset
3. For all elements whose silhouette value is within a user defined range τ calculate the EOC/ VIOC
4. Render the classified choropleth map and visually highlight all units with a VIOC value within a user defined range γ

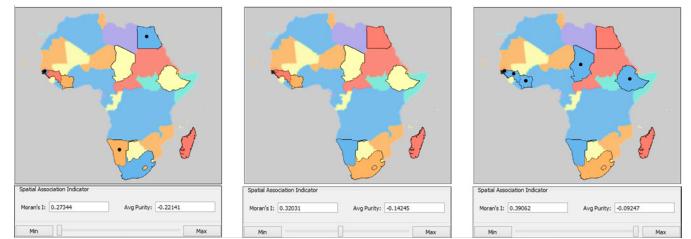


Fig. 7. Minimizing and maximizing the EOC of elements near the classification boundary using 2014 socio-economic data from the African Development Bank. 1- Minimizing the EOC, thus creating more visual heterogeneity. 2 - The initial k-means classification. 3 - Maximizing the EOC, thus creating more visual spatial clustering.

After the map is rendered, the designer can inspect the marked units, create a map that will minimize or maximize the EOC/ VIOC, manually change units near classification boundaries to obtain the desired rendering effect, or embed the EOC/ VIOC measures as uncertainty information in the map design.

4 VALIDATION

In this section, we demonstrate by example that our proposed metrics are able to identify map elements that lie near classification boundaries, whereby a small change in the boundary would impact the perceived visual spatial association.

4.1 Applying EOC for Visual Clustering

The first dataset used here is the Chicago crime data of 2014 (<https://data.cityofchicago.org>). There are 77 regions and 26 types of crime variables in this dataset. For classification, k-means clustering has been applied with $k = 5$ for three variables “Liquor Violation”, “Sex Offense” and “Robbery”. The resulting classification and the PCA scatterplot are shown in Fig. 6.2, and spatial units that may be near the cluster boundary having an impact on the EOC are highlighted by setting the silhouette value to the range of -.2 to .2. Fig. 6.1 shows the results of minimizing the EOC to reduce the spatial clustering that is visually observed, and Fig. 6.3 shows the results of maximizing the EOC to increase the spatial clustering that is visually observed.

To further summarize, our modified Moran’s I in Fig. 6.2 (the initial k-means clustering) is .14836. By shifting the classes as in Fig. 6.1, Moran’s I can be reduced to .08185 and more dispersion in the regions is seen. For example, the large purple region in the middle is dispersed as the unit on the Purple/Blue boundary shifts to Blue. Such an effect may be desirable in map design as this may help eliminate the potential for users to identify spurious patterns if the result of the visualization is designed to be as disperse as warranted by the data. Note the shift of the convex hulls in the scatterplot as well when the units are reclassified. In Fig. 6.1, we see the Purple-Blue border now overlaps as does the Red-Teal. Similarly, in Fig. 6.3, Moran’s I can be increased to .20510 and we see larger red regions form in the North and South. The Red-Teal border now overlaps in the scatterplot as well.

The most interesting boundary in the PCA projection is the Teal-Red boundary. In the Teal group, measures of the three crimes are all quite low; however, in the Red group, the data is clustered around mid-level rates. Rates are normalized by the total count of crimes, and in the Red group, liquor violations and sex offense have normalized values ranging from .19 to .43 and .21 to .47 respectively. In the Teal group these rates are 0 to .19 and 0 to .047 respectively with robbery rates in both groups being less than .21. Thus, if one were to provide a label to the Teal cluster, it could reasonably called the “low risk group” and red could be a “mid-to-high risk group”. What we see in Fig. 6 is that there are regions in the North and South of Chicago with a Teal unit surrounded by Red. When we maximize the EOC, the Red clusters become visually larger indicating more areas in the “mid-to-high risk group.” Given that the units that were changed are near the classification border, the change from Teal to Red could be warranted, and the designer’s goal could be to show that crime is a problem in

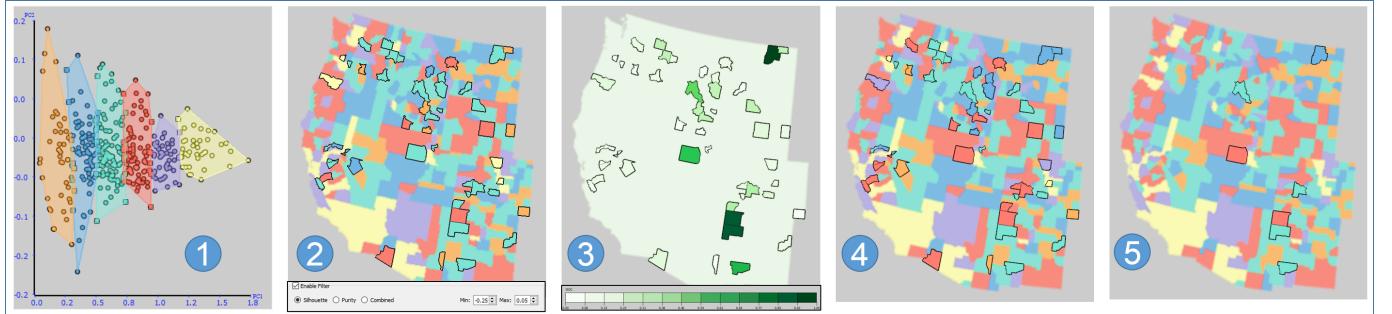


Fig. 8. Maximizing the EOC based on the VIOC near the classification boundary using US indices of industrial diversity from the western United States. 1 - The PCA scatterplot for the initial k-means classification. 2 - A choropleth map of the k-means classification, highlighted elements are those near the classification boundary. 3 - The VIOC measure of elements near the classification boundary. Darker elements will have a larger visual impact if their class changes. 4 - Maximizing the EOC of all units near the classification boundary. 5 - Maximizing the EOC of all units with VIOC in between .46 and 1.

Chicago and larger visual clusters could help sell that point. Again, the goal of this work is not on the ethical implications of such design choices, but the focus is on the fact that elements near classification borders may need to be identified to capture a holistic picture of the multivariate classification scheme. In fact, the elements that shift from Teal to Red may be some of the most interesting elements as these represent local outliers with the characteristics of more distant places. While hiding these on the map may help motivate a story, the stronger implication may be that such elements need to be highlighted to call attention to the analyst.

To further demonstrate the impact of minimizing or maximizing the EOC we explore 2014 socio-economic data from the African Development Bank Group (<http://dataportal.afdb.org/Default.aspx>). We perform k-means clustering with $k = 6$ on seven variables: the annual % of inflation; the central government's fiscal balance as a % of the GDP; the central government's total revenue and grants as a % of the GDP; the total outstanding debt as a % of the GDP; the real per capita GDP growth rate as an annual %; the gross capital formation as a % of the GDP, and; the real GDP growth as an annual %. Results of the clustering are shown in Fig. 7. Fig. 7.2 highlights all the units that can impact the EOC calculation. Fig. 7.1 is the result of minimizing the EOC. When the EOC is minimized (resulting in less visual spatial clustering), changes can be identified in Northeast Africa (the red region that was previously there has been dispersed), Northwest Africa, and the South of Africa. Fig. 7.3 is the result of maximizing the EOC. Here, more visual clustering can be observed particularly in the central blue cluster now adding contiguous members.

4.2 Combining EOC and VIOC

In Fig. 6 and Fig. 7, what becomes obvious is that the size of the spatial units plays a large role in the visual output. This is completely expected as documented in the related work [24]. Thus, while Fig. 6 and Fig. 7 focus on highlighting all units that can change with the EOC measure, the proposed VIOC measure can provide information about which units can be changed and, if changed, will have the largest visual impact. In this example, we explore measures of industrial diversity in the Western United States using data from the US Census Bureau (http://quickfacts.census.gov/qfd/download_data.html). These measures represent the relative concentration of industries for a given spatial unit of interest at a particular point in time. Fig. 8 shows the result of applying k-means clustering ($k = 6$) to the indices of health-care (N62), finance and insurance (N52), and professional and science services (N54). We use a silhouette value range of -.25 to .05 and Fig. 8.1 and Fig. 8.2 show the result of the classification with the units on a classification boundary highlighted. Note that while other units may be initially highlighted with the silhouette coefficient, by using the EOC case criteria we reduce the highlights to only those units that will have a visual impact on the map. Note that there are approximately 30 counties highlighted on the map and we want to explore which units will have the most visual impact. We use a sequential color scheme to

shade the highlighted units based on their VIOC measurement, which is directly proportional to the percent of the screen space that the spatial unit occupies. The result is shown in Fig. 8.3. As expected, the larger the county, the darker the highlighting. The reason we show this is that by simply applying silhouette filtering and EOC metrics to highlight the boundary regions, many units will be selected. If we only want to focus on the most visually salient units, the units could be further filtered based on their VIOC values. In Fig. 8.5, we set a VIOC range from .46 to 1.0 leaving only 7 of the initial 30 counties highlighted. We then modify the classification to maximize EOC. Filtering by VIOC can be thought of as another tool for the map designers' toolbox in which they can consider modifications to classes and boundaries. However, it is important to note that filtering only by size may limit the overall design space. In Fig. 8.5, we can see that in the Northeast region, the small orange county that was reclassified in Fig. 8.4 is now unchanged during the maximization of the VIOC. While the size is small, the placement creates a very strong visual break (a hole) in the cluster. As such, measures that include shape, size, distance and contiguity for filtering should also be explored as future work.

5 RECLASSIFICATION VERSUS BOUNDARY MODIFICATION

Throughout the discussion, we have primarily discussed the impact of reclassifying elements that are on classification boundaries; however, simply reclassifying an element may not be the most appropriate means of adjusting the classification. In multivariate schemes, such as k-means, recent work has focused on incorporating user feedback into the classification model [11]. Thus, if a user changes an element class, the classification model will update the weights and reassign the classification boundaries. Recent work on this topic was discussed in Sect. 2.4, and we extend our work to incorporate a modification for flexible direct manipulation.

Fig. 1 to Fig. 8 rely on what is known as result manipulation, which means the modification will only effect on the class index of the user selected units. Model manipulation, on the other hand, will affect the weights in the clustering processing and eventually the classes of other data points. Each element in the dataset will have an associated weight that can be modified through user interaction. Suppose there are n units u_1, u_2, \dots, u_n and their weights are formed as w_1, w_2, \dots, w_n , such that initially each spatial unit u_i will have the same instance weight $w_i = 1$ influencing the placement of the centroids. After the initial clustering, analysts may assign unit u_i to specific cluster C_j , then u_i 's weight will be modified to be either based on the cluster C_j 's size s_j such as $w_i = 1 + \sqrt{s_j}$ or a predefined constant value larger than 1. During each iteration, u_i 's proximity to C_j 's centroid c_j is increased by multiplying that weight w_i . Thus u_i is more likely to be assigned to cluster C_j . When calculating the cluster centroid, u_i will only contribute its weight w_i to the cluster it belongs to. The new centroid of C_j will be $c_j = \frac{\sum_{i \in J} w_i u_i}{\sum_{i \in J} w_i}$ where J is the set of units that have been assigned to C_j . Eventually, this result in cluster C_j 's centroid c_j moves towards u_i and this unit will likely belong to that cluster.

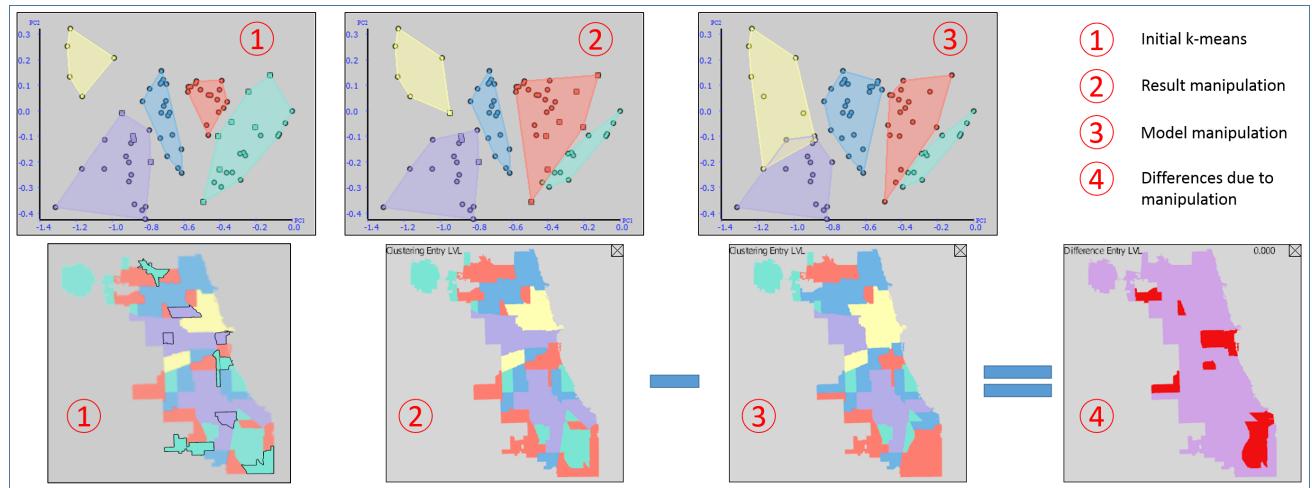


Fig. 9. The effects of model manipulation on choropleth map classification. 1 - A k-means classification of criminal incident reports in Chicago, IL. 2 - Maximizing the EOC through result manipulation (i.e., changing the unit classification does not affect the k-means weight). 3 - Maximizing the EOC through model manipulation (i.e., changing the unit classification updates the k-means weights). 4 - The difference between result manipulation and model manipulation. Note that units that were not originally highlighted as being near the classification boundary are now reclassified due to the updated weights used in k-means.

We can also apply a model manipulation scheme. First, we identify elements to reclassify using the silhouette range τ . We can then maximize (or minimize) the EOC that forces the reclassifying of elements. This reclassifying will automatically update the weights of the k-means clustering, and a new classification based on the updated weights will be generated. This result is shown in Fig. 9. Here, we revisit the data and classification scheme applied in Fig. 6 (the Chicago crime data).

Fig. 9.1 and Fig. 9.2 are the same k-means and maximized EOC results from Fig. 6.1 and Fig. 6.4 respectively. What is interesting is that by changing the k-means weights, the classification boundaries shift and units that were not marked as boundary candidates are now subsumed by a new class. In Fig. 9.2, the same units that were reclassified in Fig. 9.2 are reclassified in Fig. 9.3. This causes the weights in the k-means clustering to update, and then a new k-means classification is performed, resulting in the map classification of Fig. 9.3. If we take a difference between Fig. 9.2 and Fig. 9.3, we can see what other units were shifted as a result of updating the weights of the k-means classification (Fig. 9.4), and we notice that an even larger amount of visual spatial clustering can be seen in Fig. 9.3 than in Fig. 9.2.

6 CONCLUSION AND FUTURE WORK

As previously stated, a critical step in designing choropleth maps is the choice of classification method. How a map is classified directly impacts the resultant visual output and can lead to misinformation about the underlying data. In order to assess the visual impact of such choices, we have developed a methodology for quantifying the visual impact of adjusting classification boundaries in a choropleth map and present a scheme for maximizing or minimizing the amount of visual clustering present in the map and demonstrated the results using several datasets. What is important to note is that the goal of choosing classification boundaries is to achieve a reasonable split in the data, and this is often left up to the designer. By providing designers with new ways to assess the visual impact of small classification changes, the designer can further refine and assess their map message.

Perhaps the most intriguing part of our results is the output when applying a reweighting scheme to the classification. While the reweighting of units in Fig. 9 resulted in more changes (and arguably more visual clustering) than a naïve reclassifying of elements, it is likely that this reweighting could also introduce a reduction of visual clustering. Future research should explore the sensitivity of such an application across various clustering schemes. Furthermore, research into what pattern changes result in an increased perception of visual clustering should also be undertaken. In this paper, we rely solely on the fact that colors are changing and regions are becoming larger. Past research [24] has

shown that the larger a patch of color becomes in a choropleth map, the more likely that it will be identified as a cluster. However, there may be particular patches that could be changed that may have a greater impact on the perception. Of course, the size of the spatial unit matters a great deal, but what if changing the classification of a spatial unit fills in a donut hole? Is this perceived as resulting in more spatial clustering than if we change a spatial unit's classification such that it just adds to the edge of the donut? What if a spatial unit acts as a bridge? For example, if there are two spatial groupings with the same class separated by a narrow band of other classes, how is the spatial clustering perceived if one unit is changed to create a bridge? Understanding the impact of these patterns would allow us to computationally identify them and use these types of patterns to create a more perceptually rigorous VIOC metric. Future work will focus on the use of such metrics for highlighting uncertainty within the map, as well as exploring boundary elements with respect to statistical measures of spatial clustering. Specifically, if a region is found to be a statistically significantly spatially cluster, should boundary elements contiguous to this region be adjusted to highlight the significance?

While our study presents an initial methodology, there is a wide range of future research that needs to be undertaken to further explore the visual impact of classification boundaries in choropleth maps. First, the class stability metric applied here was designed to be relatively agnostic to the underlying clustering algorithm used; however, different clustering algorithms have different properties and constraints. For example, in the k-means approach showcased here, the choice of seed points may greatly impact the boundary elements, as would the choice of single-linkage versus complete-linkage in hierarchical clustering. It is likely that clustering specific methods might better capture the underlying uncertainties associated with the resulting visual representation. Second, our VIOC and EOC measures focused only on neighboring spatial associations and size. Future work could explore the resultant shape of a cluster or the area of an element as well as metrics that take spatial distances into account too. Finally, this research points to the need for a comprehensive study on this issue where the underlying data distributions vary along with the proposed clustering algorithm in order to generate comprehensive guidelines for how to approach boundary issues given a known data distribution and clustering methodology.

ACKNOWLEDGMENTS

We would like to thank the reviewers for their suggestions that greatly improved the manuscript as well as Michael Steptoe for his work in video and demo production. This work was supported by the NSF under Grant No. 1350573.

REFERENCES

- [1] G. Andrienko and N. Andrienko. Interactive maps for visual data exploration. *International Journal of Geographical Information Science*, 13(4):355–374, 1999.
- [2] G. Andrienko, N. Andrienko, S. Rinzivillo, M. Nanni, D. Pedreschi, and F. Giannotti. Interactive visual clustering of large collections of trajectories. In *IEEE Symposium on Visual Analytics Science and Technology*, pp. 3–10. IEEE, 2009.
- [3] G. Andrienko, N. Andrienko, and A. Savinov. Choropleth maps: classification revisited. In *Proceedings of ICC*, pp. 6–10, 2001.
- [4] L. Anselin. Local indicators of spatial association - LISA. *Geographical Analysis*, 27(2):93–115, 1995.
- [5] M. P. Armstrong, N. Xiao, and D. A. Bennett. Using genetic algorithms to create multicriteria class intervals for choropleth maps. *Annals of the Association of American Geographers*, 93(3):595–623, 2003.
- [6] D. I. Ashby and P. A. Longley. Geocomputation, geodemographics and resource allocation for local policing. *Transactions in GIS*, 9(1):53–72, 2005.
- [7] H. G. Basara and M. Yuan. Community health assessment using self-organizing maps and geographic information systems. *International journal of health geographics*, 7(1):1, 2008.
- [8] A. Ben-Hur, D. Horn, H. T. Siegelmann, and V. Vapnik. Support vector clustering. *The Journal of Machine Learning Research*, 2:125–137, 2002.
- [9] B. Boots. Developing local measures of spatial association for categorical data. *Journal of Geographical Systems*, 5(2):139–160, 2003.
- [10] C. A. Brewer and L. Pickle. Evaluation of methods for classifying epidemiological data on choropleth maps in series. *Annals of the Association of American Geographers*, 92(4):662–681, 2002.
- [11] K. Chen and L. Liu. VISTA: Validating and refining clusters via visualization. *Information Visualization*, 3(4):257–270, 2004.
- [12] K. Chen and L. Liu. iVIBRATE: Interactive visualization-based framework for clustering large datasets. *ACM Transactions on Information Systems (TOIS)*, 24(2):245–294, 2006.
- [13] J. Cheshire, P. Mateos, and P. A. Longley. Delineating europe’s cultural regions: Population structure and surname clustering. *Human Biology*, 83(5):573–598, 2011.
- [14] A. D. Cliff and K. Ord. Spatial autocorrelation: A review of existing and new measures with applications. *Economic Geography*, 46:269–292, 1970.
- [15] E. Cromley and R. Cromley. An analysis of alternative classification schemes for medical atlas mapping. *European Journal of Cancer*, 32(9):1551–1559, 1996.
- [16] M. F. Dacey. *A review on measures of contiguity for two and k-color maps*. Dept. of Geography, Northwestern University, Evanston, Illinois, 1965.
- [17] S. L. Egbert and T. A. Slocum. EXPLOREMAP: An exploration system for choropleth maps. *Annals of the Association of American Geographers*, 82(2):275–288, 1992.
- [18] I. S. Evans. The selection of class intervals. *Transactions of the Institute of British Geographers*, pp. 98–124, 1977.
- [19] R. Geary. The contiguity ratio and statistical mapping. *The Incorporated statistician*, 5(3):115–146, 1954.
- [20] A. Getis and J. K. Ord. The analysis of spatial association by use of distance statistics. *Geographical Analysis*, 24(3):189–206, 1992.
- [21] K. Goldsberry and S. Battersby. Issues of change detection in animated choropleth maps. *Cartographica: The International Journal for Geographic Information and Geovisualization*, 44(3):201–215, 2009.
- [22] D. Guo. Regionalization with dynamically constrained agglomerative clustering and partitioning (redcap). *International Journal of Geographical Information Science*, 22(7):801–823, 2008.
- [23] D. Guo, J. Chen, A. M. MacEachren, and K. Liao. A visualization system for space-time and multivariate patterns (VIS-STAMP). *IEEE Transactions on Visualization and Computer Graphics*, 12(6):1461–1474, Nov. 2006.
- [24] M. Haklay. *Interacting with geospatial technologies*. Wiley Online Library, 2010.
- [25] W. W. Hargrove and F. M. Hoffman. Potential of multivariate quantitative methods for delineation and visualization of ecoregions. *Environmental Management*, 34:S39–S60, 2005.
- [26] S. Haroz and D. Whitney. How capacity limits of attention influence information visualization effectiveness. *IEEE Transactions on Visualization and Computer Graphics*, 18(12):2402–2410, 2012.
- [27] R. Harris, P. Sleight, and R. Webber. *Geodemographics, GIS and neighbourhood targeting*, vol. 7. John Wiley and Sons, 2005.
- [28] J. A. Hartigan and M. A. Wong. Algorithm AS 136: A k-means clustering algorithm. *Applied statistics*, pp. 100–108, 1979.
- [29] L. J. Hubert, R. G. Golledge, and C. M. Costanzo. Generalized procedures for evaluating spatial autocorrelation. *Geographical Analysis*, 13(3):224–233, 1981.
- [30] G. F. Jenks. The data model concept in statistical mapping. *International yearbook of cartography*, 7(1):186–190, 1967.
- [31] A. Klippen, F. Hardisty, and R. Li. Interpreting spatial patterns: An inquiry into formal and cognitive aspects of tobler’s first law of geography. *Annals of the Association of American Geographers*, 101(5):1011–1031, 2011.
- [32] J. Krygier and D. Wood. *Making maps: A visual guide to map design for GIS*. Guilford Press, 2011.
- [33] R. Lloyd and T. Steinke. Visual and statistical comparison of choropleth maps. *Annals of the Association of American Geographers*, 67(3):429–436, 1977.
- [34] R. E. Lloyd and T. Steinke. The decisionmaking process for judging the similarity of choropleth maps. *The American Cartographer*, 3(2):177–184, 1976.
- [35] P. Longley. *Geographic information systems and science*. John Wiley & Sons, 2005.
- [36] P. A. Longley. Geodemographics and the practices of geographic information science. *International Journal of Geographical Information Science*, 26(12):2227–2237, 2012.
- [37] E. M. Macambira and C. C. de Souza. The edge-weighted clique problem: Valid inequalities, facets and polyhedral computations. *European Journal of Operational Research*, 123(2):346–371, 2000.
- [38] S. Moon, E.-K. Kim, and C.-S. Hwang. Effects of spatial distribution on change detection in animated choropleth maps. *Journal of the Korean Society of Surveying, Geodesy, Photogrammetry and Cartography*, 32(6):571–580, 2014.
- [39] P. A. Moran. Notes on continuous stochastic phenomena. *Biometrika*, pp. 17–23, 1950.
- [40] T. Nakaya, A. S. Fotheringham, C. Brunsdon, and M. Charlton. Geographically weighted poisson regression for disease association mapping. *Statistics in medicine*, 24(17):2695–2717, 2005.
- [41] J. Olson. *The Effects of Class Interval Systems on the Visual Correlation of Choropleth Maps*. PhD thesis, University of Wisconsin–Madison, 1970.
- [42] J. M. Olson. Autocorrelation and visual map complexity. *Annals of the Association of American Geographers*, 65(2):189–204, 1975.
- [43] M. Polczynski and M. Polczynski. Using the k-means clustering algorithm to classify features for choropleth maps. *Cartographica: The International Journal for Geographic Information and Geovisualization*, 49(1):69–75, 2014.
- [44] P. J. Rousseeuw. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53–65, 1987.
- [45] M. W. Scripter. Nested-means map classes for statistical maps. *Annals of the Association of American Geographers*, 60(2):385–392, 1970.
- [46] A. D. Singleton and P. A. Longley. Geodemographics, visualisation, and social networks in applied geography. *Applied Geography*, 29(3):289–298, 2009.
- [47] A. Slingsby, J. Dykes, and J. Wood. Exploring uncertainty in geodemographics with interactive graphics. *IEEE Transactions on Visualization and Computer Graphics*, 17(12):2545–2554, 2011.
- [48] M. Sun, D. Wong, and B. Kronenfeld. A heuristic multi-criteria classification approach incorporating data quality information for choropleth mapping. *Cartography and Geographic Information Science*, pp. 1–13, 2016.
- [49] D. Vickers and P. Rees. Creating the UK National Statistics 2001 output area classification. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 170(2):379–403, 2007.
- [50] C. Xia, W. Hsu, M. L. Lee, and B. C. Ooi. Border: Efficient computation of boundary points. *IEEE Transactions on Knowledge and Data Engineering*, 18(3):289–303, 2006.
- [51] N. Xiao and M. P. Armstrong. Supporting the comparison of choropleth maps using an evolutionary algorithm. *Cartography and Geographic Information Science*, 32(4):347–358, 2005.
- [52] Y. Zhang, W. Luo, E. A. Mack, and R. Maciejewski. Visualizing the impact of geographical variations on multivariate clustering. *Computer Graphics Forum*, 35(3):101–110, 2016.