

Social Area Analysis, Data Mining, and GIS

Seth E. Spielman, Department of Geography, SUNY-Buffalo, Buffalo, NY,
ses27@buffalo.edu

and

Jean-Claude Thill, Department of Geography and Earth Sciences, University of North
Carolina at Charlotte, Charlotte, NC, jfthill@uncc.edu

Abstract:

There is a long tradition of describing cities through a focus on the characteristics of their residents. A brief review of the history of this approach to describing cities highlights some persistent challenges. To fit within the constraints of widely used multivariate data reduction techniques and thematic cartography populations are classified. Historically, this classification has been guided by theory (i.e. Shevky-Bell) or simply the desire to efficiently describe urban populations. The labeling of classes reduces the very complexity these multivariate methods and maps are trying to capture. This problem is addressed through a geodemographic approach that uses the Kohonen Self-Organizing Map (SOM) algorithm. A dataset describing 79 attributes of the 2217 census tracts in New York City is analyzed through a method that pairs a SOM with a GIS. The resulting “maps” represent social space and geographic space and can be used to explore the similarities and differences among regions of the city and the geographic patterns of different geodemographic classes. Pairs of social and geographic maps are formally compared using simple pattern metrics.

Key Words: Self-Organizing Maps, Geodemographics, New York City, Data Mining, GIS

Introduction

In gearing up for the first United States decennial census in 1790, James Madison argued that the census should be "extended so as to embrace some other objects besides the bare enumeration of the inhabitants; it would enable them to adapt the public measures to the particular circumstances of the community" (Kurland & Lerner, 1987). Madison's idea, that knowing something about the characteristics of local populations improves local governance is accepted as a basic premise in planning, politics, and policy analysis. However how one understands the particular circumstances of a community is a methodological question that has been evolving for over a century.

Madison's proposal to extend the census to include the occupations of inhabitants was rejected by the United States Senate in 1790. In a letter to Jefferson, Madison reflected that his plan was "thrown out by the Senate as a waste of trouble and supplying materials for idle people to make a book" (Cohen, 1981). Unlike in Madison's day, data about cities and the people who live in them is now abundant; in fact data are so abundant and complex that integrating available information into the public planning processes is often difficult. The first census asked five questions; the long form of the questionnaire for the 2000 decennial census of population was 10 pages long and included over 50 questions. Many municipalities now maintain detailed datasets describing crime, traffic, school performance, the built environment, and many other facets of urban life. The volume of data currently available to planners is excellent fodder for urban scholars. Yet, it still remains a great challenge to communicate the complexity of the urban social landscape to the public in an engaging and efficient manner.

In addition to a dramatic increase in the volume of information, new forms of analysis that emphasize an exploratory approach and are based on computational principles have become commonplace. Data mining is "the extraction of implicit, previously unknown, and potentially useful information from data" (Witten & Frank, 2005 p. xxiii). Machine learning techniques of data mining, while still seldom used in urban analysis, have the potential to help analysts develop detailed differentiation of the urban landscape. In contrast to more conventional multivariate statistical methods such as factor analysis, principal component analysis, and multidimensional scaling, they tend to be less bound by a priori assumptions. On the other hand, Geographic Information Systems have become an indispensable tool of urban analysis in large part because they facilitate cartographic visualization and management of geographically referenced data. This paper presents a novel application of geographic information systems by integrating them with data mining to characterize populations in urban areas using large datasets. The Kohonen Self-Organizing Map algorithm (Kohonen, 1997) is applied to a dataset containing 79 attributes describing census tracts in New York City. The result is a typology of census tracts presented as a pair of linked maps- one representing social space and another representing geographic space. These maps capture the complexity of New York's social landscape and provide insight into the relationship between geographic proximity and social similarity at the tract level.

Our goal in this paper is to present a new approach to the problem of describing communities through a focus on the characteristics of residents. The history of residential segregation by race and income in America has supported the use of very general colloquial descriptions of neighborhoods that focus almost exclusively on combinations of these two factors. It is our contention in this paper that using large amounts of data to describe

neighborhoods, while it presents technical challenges, has the potential to improve planning and governance through a better understanding of the “particular circumstances of the community.” We are not the first to make this argument. We will review the history of efforts to describe neighborhood demographics, contrast our approach with the prevailing multivariate methods (geodemographics), and present our method and results.

Maps and Neighborhood Typologies

Since the turn of the previous century, advocates and social scientists have been mapping the socio-economic variation in cities through looking at residential patterns. Charles Booth’s poverty maps of London are a classic effort to map this social landscape. Booth, working between 1886 and 1903, classified London’s streets as using seven categories: wealthy, well-to-do, fairly comfortable, mixed, poor, very poor, and vicious, semi-criminal (Booth, 1902).

In spite of the abundance and complexity of spatial data describing the US population in the planning and policy context, one often finds that we have not moved very far beyond Booth’s classification system. Neighborhoods are often differentiated using just a few attributes- the income, race, and occupation of inhabitants and the density of the built environment. Descriptive terms like “working class suburbs” and “poor inner city” evoke images of prototypical neighborhoods. Among the residents of a given city, neighborhood names are often signifiers of subtle differentiations in social and physical landscape. These subtle distinctions are often hard to communicate to non-residents and may not be commonly understood by residents.

In the modern context, the most sophisticated efforts to classify populations are known as geodemographic or market segmentation systems. Geodemographic systems classify small areas into discrete categories using consumer behavior, lifestyle, and demographic data. These tools are widely used in the commercial sector and multivariate social classifications of “neighborhoods” has become an international industry (Harris et al., 2005; Longley & Clarke, 1995). Unfortunately, since market segmentation is a competitive, commercial enterprise, the specifics of the methods and data used in the construction of these proprietary systems is often obscure.

As described by Harris et al. (2005), the geodemographic approach is essentially a data reduction technique. A standard methodology involves a double weighted k-means algorithm to develop initial clusters. K-means is a simple type of cluster analysis where the user chooses a desired number of clusters, k, and then observations are assigned to clusters based on their proximity to the cluster means. The initial cluster centers can be randomly assigned by the algorithm or manually specified a priori. The procedure is iterative and generally ends once the clusters become stable. In the approach outlined by Harris et al., areas are weighted based on their population and variables are weighted based upon how important the variable is in distinguishing different types of consumers. The result is that neighborhoods (usually treated as zip codes or census tracts) are assigned to a predetermined number of clusters representing similar types of neighborhoods.

Geodemographers then assign evocative titles to each cluster. Names like “American Dreams” and “Multi-Culti Mosaic” are used by the Prizm lifestyle segmentation system (Weiss, 1989; Curry, 1993). These clusters are described with “pen pictures,” which are short 1-paragraph descriptions of the discriminating characteristics of each cluster. The user

interacts with the system through the category titles. Geodemographic systems divide the national population into discrete classes on the basis of variables useful for describing consumption patterns; they are tools primarily designed for “differentiating between different categories of rich people” (Webber, 2004, p. 220). Whereas, Charles Booth was interested in differentiating different classes of poor people (worthy vs. unworthy poor) (Ward, 1990). A balanced classification of the city’s residents would treat all groups equally instead of segmenting some populations more than others.

Batey and Brown (1995) and Harris et al. (2005) see modern geodemographic systems as rooted in a conceptualization of neighborhood based on the human ecologic approach of the Chicago School: “geographical units distinguished by both physical individuality and by the social and cultural characteristics of the population” (Batey & Brown, 1995). The approach used by Booth and geodemographic systems, i.e., named categories, is not the only way to develop multi-dimensional classifications of urban areas. Some of the earliest classifications of census tracts in the human ecologic tradition were done by Eshref Shevky in the 1940s in Los Angeles using seven variables and over three hundred census tracts. Shevky created three indices by computing percentiles for seven variables. The first index measured urbanization, the second measured segregation, and the third measured “social rank.” Shevky ranks places and then compares places to each other and to the city-wide average. He also groups places that have similar ranking on each of the three dimensions. His 1949 book includes extensive tables reporting these results. Interestingly, these classifications were a response to what Shevky saw as the trend in urban studies toward descriptive generalizations. Shevky hoped that by developing a typology of urban places through a focus on local characteristics one could build a more robust understanding of urban systems (Shevky & Williams, 1972).

Shevky's early work on social area analysis was instrumental in the emergence of “factorial ecology” as a line of inquiry. The term, factorial ecology, emerged in the mid-60s and refers to the use of factor analysis to differentiate areal (ecological) units using the characteristics of residents (Janson, 1980). Factorial ecologies most typically describe the characteristics of urban areas through an analysis of census tracts, however, during the heyday of factorial ecology, factor analytic approaches were used to describe patterns of areal differentiation at various geographic scales (Berry, 1971; Rees, 1971). Hundreds of studies have been published in the factor ecologic tradition. As an indication of the popularity of the approach, the bibliography of a 1971 special issue of *Economic Geography* on comparative factorial ecology included over 350 references.

Factor analysis is a method to reduce a large matrix of units of observation and their attributes to a smaller number of factors. Berry (1971) uses the following analogy to describe the factor analytic approach (p. 215-216):

“If there are n areas and m variables, an $n \times m$ matrix is used to list the manifest evidence. An atlas comprising m plates could also depict the variations. Factorial methods are brought into play to determine the latent structure of dimensions of variation - the repetitive sequences - underlying the manifest experiences of the atlas.”

The smaller matrix is a more concise description of the economic and demographic variability of census tracts. Factor scores are sometimes described as “latent” or “fundamental” variables. Interpretation of latent variables is a matter of some debate, some

use latent variables to explain urban residential patterns (Ward, 1969) while others simply saw them as concise descriptions of patterns (Rees, 1971, p. 221). The former view is particularly controversial.

In the explanatory mode factor scores are interpreted as representations of theoretical constructs (Berry, 1971; Janson, 1980). Theory held that residential areas in Western industrialized cities, particularly those in the United States, were differentiated by three factors; one describing racial and ethnic segregation, another describing socioeconomic status, and a third describing one's point in the lifecycle. This three-factor view is rooted in Shevky's early analysis of Los Angeles and is known as the Shevky-Bell hypothesis (Janson, 1980). While factor analysis is not a confirmatory statistical technique, the fact that some form of the Shevky-Bell factor structure emerged from many urban analyses was seen as support for this view of urban spatial structure. Palm and Caruso (1972) saw this argument as a form of "speculative synthesis." A factor consists of many variables, each one weighted differently. Palm and Caruso argue that the labels used to describe factor scores generally focus on only a few of the variables loaded on that particular factor (Palm & Caruso, 1972). Their indictment of factor analysis is extensive and beyond the scope of this paper. For our purposes here, it is interesting to note that their criticism of the "crudeness of classification" in factor analysis could be extended both to modern geodemographic systems and early geographic studies of urban populations.

The goal of commercial geodemographic packages is to place local areas in some national context based on the characteristics of residents, that is, their primary purpose is descriptive generalization. As described earlier, this purpose is satisfied by labeled categories. As national classifications have proven useful for marketing and are widely used as predictors of consumer behavior (Webber, 1985), the authors do not wish to challenge the utility of geodemographic systems in marketing research. However, when one looks at the history of efforts to map the socioeconomic variation in cities certain themes emerge. Labeled categories have been used for over a hundred years to describe urban populations in a multivariate sense. The act of labeling may be problematic and it may reduce or oversimplify the very complexity that these multivariate mapping efforts seek to capture. While the techniques have evolved and become more sophisticated, while the volume and perhaps quality of the data has greatly increased, the basic method of multivariate mapping has not changed. For as long as such maps have been made, labeled categories have been used.

Our thesis is that the inductive, quantitative analytical techniques that have been used since Shevky do not require the use of labeled categories. The only important distinctions between this work and earlier efforts at describing urban populations is that the techniques employed here allow one to avoid the use of labeled categories, assess the multivariate similarity of classes, and explore the relationship between proximity and similarity. The goal of this paper is not to develop a new theory of urban residential spatial structure nor is it to examine the process underlying residential aeral differentiation. Rather it is an effort to apply modern methods of data exploration and analysis to the very old problem of describing urban populations.

Self-Organizing Maps

Maps preserve topological relationships among objects in space. In the cartographic context, entities and features that are close to each other in the real world are represented close to

each other on a map. There is evidence to suggest that the ability to situate oneself on a map is an innate human ability (Holden, 2006). This makes maps useful tools for describing the environment and presenting data. Maps are frequently used to present information about urban areas. Traditional cartographic maps are limited in that they can only paint a one-dimensional picture of the social characteristics of an area. While maps are an efficient and familiar medium, they have limitations when it comes to displaying multiple pieces of information about the same location.

The concept of a map can also be applied to non-geographic objects; or it can be used to visualize geographic objects (census tracts) in a spatial but non-geographic context. That is, census tracts can be organized in space based upon the similarity of their characteristics rather than their geographic proximity. This is the basic idea behind the Kohonen algorithm that creates Self-Organizing Maps (SOM). In the early 1980s, Teuvo Kohonen developed a technique to map similar patterns onto contiguous areas in output space. The resulting visualizations are called self-organizing feature maps (Kohonen, 1997). The idea is simple: observations (vectors) that are similar are mapped to proximate regions of a two-dimensional synthetic space of fixed topology. SOMs are a type of unsupervised artificial neural network. Neural networks use the concept of a “neuron” to analyze data. Neurons are organized in layers and connected. Neurons respond to a stimulus (data) by transforming the data, themselves, or other neurons. In the approach outlined by Kohonen (1997) and used here, a single output layer of neurons is trained such that regions of this layer are sensitized to observations with specific types of attribute vectors.

SOM outputs are attribute maps. Unlike thematic maps, SOM feature maps excel at the display of high dimensional datasets. Feature maps are a projection of high dimensional attribute space such that attribute vectors of a particular generalized form are associated with locations in output space (Skupin & Agarwal, 2004; Skupin & Fabrikant, 2003). As a data reduction method, a SOM cuts down the number of rows and columns of a data matrix; the method is a combination of data projection and data quantization (Yan & Thill, 2007). With a self-organizing feature map, a map-reader can judge the similarity or dissimilarity of objects based on their proximity. The approach shares some characteristics with multidimensional scaling, regression, and cluster analysis. The process of fitting observations to a SOM is an iterative and stochastic process dependent upon a random map initialization. For details on specifying and training a SOM see (Bacao et al., 2001; Kohonen, 1997; Kohonen et al., 1995), Openshaw (1989), and others. This paper will only describe the details relevant to the interpretation of SOMs.

Space in a SOM consists of a regular lattice of “neurons” each of which stores a vector describing attribute weights. The elements of the lattice generally are square or hexagonal. Through the SOM mapping process, each neuron in the output layer is sensitized to a particular configuration of attributes and observations are “fit” to neurons much as a regression model is fit to data. It is useful to think of the neurons on the feature map as buckets for data. Observations that are similar are placed either in the same bucket or in buckets that are topologically close to one another on the feature map. For example, places with many wealthy householders, with high levels of education, high homeownership rates, and low poverty rates would end up in buckets that are near each other and clustered in a region of the feature map. On the other hand, census tracts where poverty is abundant and residents typically have low levels educations would end up clustered in buckets in a different

region of the SOM; probably quite far away from the well educated and wealthy people. Places that have both wealthy households and poor households would end up occupying a region of the map somewhere between the two extremes. Training a SOM is an iterative process of defining what types of observations are associated with buckets in different regions of the feature map. By examining the contents of each bucket after the SOM is completely trained, one can get a sense of how different regions of the SOM represent different types of observations.

SOM feature maps of different sizes have different characteristics (Skupin & Agarwal, 2004). Small feature maps provide generalizations; large grids allow a unique location in geographic space to be mapped to a unique location in the synthetic attribute space. In a large feature map, where the number of buckets exceeds the number of observations, each bucket may hold few, if any observations; regions have very specific properties. On the contrary, in a small SOM feature map where the number of observations far exceeds the number of buckets, many observations will fall into each bucket and regions of the map will represent general characteristics (FIGURE 1). The size of the SOM feature map is specified by the user a priori.

FIGURE 1 ABOUT HERE

Relatively few geographic applications of SOM have so far been reported in the literature. The SOM has successfully been trained to classify digital satellite images (Villmann et al., 2003; and many others). In all these works, SOM is used as an unsupervised classifier, working on the multi-spectral information in satellite images. Openshaw and Wymer (1995) tested an application of the SOM algorithm against a K-means classification on census data in the United Kingdom. Skupin and Hagelman (2003) study patterns of change in the socioeconomic profile of neighborhoods with the SOM method applied to census tract statistics. Outside the application of SOM to satellite imagery or census data, a handful of studies of geographic feature identification have been conducted with the SOM method. An early case study by Kaski and Kohonen (1996) applied SOM to a data set of 39 welfare statistical indicators of countries. Himanen et al. (1998) explored the applicability of SOM in identifying daily travel patterns in a disaggregate travel diary data set. Thill et al. (2007) analyze ill-conditioned linguistic data on the Atlantic Seaboard of North America in relation to geography. Yan and Thill (2007) developed an interactive visual data mining environment to explore patterns in a multidimensional database of air travel flows. Kauko (2005) studied spatial housing markets, and Hatzichristos (2004) applied SOM to a regional classification of Athens, Greece.

Skupin and Hagelman (2003, 2005) used large grids to explore the demographic “trajectories” of different regions of Texas. In this context, a large grid separates similar regions into unique areas on the feature map. In this work, a SOM is trained on 30 years of census data. The large feature map allowed them to examine how the characteristics of census tracts changed over time by looking at how individual tracts moved around the output space over time. Medium sized grids are a compromise; they allow regions with clearly identifiable characteristics to form on the map, yet general statements can be made about these regions as they contain a fair number of census tracts— this is the approach we used for our case study, which we will discuss in the next section.

The Kohonen Self Organizing Map algorithm extends modern geodemographics, and similar cluster-like methods by constructing topological relationships between classes (Kohonen, 1997). Geodemographic classifications group areas with similar characteristics and apply descriptive labels to these classifications. One of the problems with such classifications is that groups are discrete. It is not clear how similar or dissimilar classes are in a multivariate sense because classes are typically described by comparison to regional or national averages using income and consumer behavior. By constructing topological relationships between classes, the Kohonen algorithm allows the user to understand the degree of similarity or dissimilarity between areas based upon their location in a two dimensional projection of multidimensional attribute space. With an integrated visual data mining approach, we avoid the use of category labels. Since our approach is visual, we can define a very large number of categories, over 1000, and still present our results in a way that is easy to interpret.

To explore the efficacy of SOM and geovisualization as a geodemographic tool, a dataset with 79 variables is used to describe census tracts in New York City. The variables used in the rest of the analysis are listed in the Appendix. Variables are selected from all 2000 Population Census variables to represent some aspect of New York's social geography. Variables are mapped onto a 45 x 30 output map consisting of 1350 buckets (neurons) for 2217 census tracts. Buckets can be interpreted as classes or clusters of similar data.

Mapping New York City

New York City is an ideal subject for testing for spatial demographic methods. New York is home to what may be the most racially and ethnically diverse zip code in the nation, 11373 in Elmhurst (a neighborhood in Queens County) where the local high school has students from 96 different nations and that speak 59 languages (Utley, March 17, 2001). New York also has clearly identifiable ethnic enclaves. New York has well defined high-income areas. Some of the wealthiest parts of the United States are in the city, yet the Bronx is the poorest urban county in the nation. This combination of diversity and residential segregation make simple low dimensional classifications of New York's neighborhoods difficult. The complexity and richness of New York's social landscape make it ideally suited to exploration through data mining techniques and geovisualization tools.

The SOMPAK code library was used and a SOM was trained using random selection of 50% of the census tracts (Kohonen et al., 1995). Parameterization of each step has a large effect on the resultant trained map. Training a SOM is more akin to an art than to a science, hence the widely held view that SOMs, like other data mining techniques, are "black boxes" (Miller & Han, 2001). We chose suitable SOM parameters through trial and error. The final map was selected through an iterative process whereby we initialized 100 SOMs using random numbers and trained each SOM by presenting 100,000 census tracts (the training dataset was sampled with replacement). The map with the lowest mean square error was retained for analysis. This training period is computationally intensive and took 8 run-time hours on a desktop computer with an AMD Athalon XP 3200 processor and 1GB of RAM. The SOM output was imported into ESRI ArcGIS 9.1 software using a Python script. Using the ESRI geodatabase file format, a relational (one to many) link was established between the self-organizing feature map and a geographic map of New York City by census tracts.

FIGURE 2 ABOUT HERE

There are a number of different ways to summarize a SOM. Traditionally, component planes and the unified distance matrix (or U-matrix) are utilized. The U-matrix is a visualization of the SOM that illustrates the distance between adjacent neurons in attribute space (the U-matrix for the SOM described below is shown in Figure 2). Observations that have similar profiles on input variables are mapped to nearby areas, however, distance in the synthetic space of the SOM is not constant. Some pairs of proximate buckets may hold observations that are more similar than other pairs of proximate buckets. The synthetic space of the SOM has hills and troughs which can increase surface distance between pairs of proximate neurons. The U-matrix shows the distance, or dissimilarity, between the vectors describing each neuron, it illustrates cluster structures evidenced by “troughs” and “hills” in the distance surface. The U-matrix provides little insight on the meaning of the observed structures. Each bucket in the SOM has a unique value for each of the 79 attributes in the data set. A component plane uses color to represent the weight assigned to a single input variable at each neuron. Therefore, inspection of each of the 79 component planes would in principle allow a user to figure out the exact characteristics of each neuron. This approach however holds little advantage over an atlas displaying the same data.

An alternative approach is to work backwards, that is by selecting a census tract or group of tracts with known characteristics and examining where they fall on the SOM feature map. In applied settings the relative difference and/or similarities between census tracts is often of interest. One can then use knowledge of the city under study to explore the geography of the SOM feature map. By selecting an area of interest one can examine how it maps onto the SOM feature map. Reversing the process, selecting all tracts that fall into the same buckets as the area of interest, lets one quickly visualize parts of the city that are similar to the area of interest in a multivariate sense. Figure 3 illustrates the latter approach to SOM-based urban social geography. The census tracts in Manhattan’s Community Board 8 (an administrative unit that has a role in governance) are selected in this figure. The 32 tracts that make up Community Board 8 map to a relatively well-defined region of the SOM feature map illustrating that Community Board 8 is a (relatively) socially homogenous political unit. Community Board 8 is one of the most affluent in the city encompassing areas to the east of the southern half of Central Park in Manhattan. Most of the 32 tracts are mapped to 26 neurons bundled together in the upper right region of the SOM feature map. However, two or three outliers are visible. Those outliers correspond to census tracts in Community Board 8 that contain public housing developments. The rightmost of the three outliers is a tract that contains a development for low income senior citizens. The two furthest outliers each contain large, high-rise, low income housing projects (Isaacs Towers and John Haynes Holmes Towers). Given their discordant socio-economic profiles, these tracts sensitize distant parts of the SOM feature map, in spite of their close geographic proximity to the rest of Community Board 8. The third image in the sequence illustrates the return to the geographic map, where twelve additional census tracts that are similar, i.e. occupying the same part of the SOM, to those in Community Board 8 are identified. Most of these new tracts are geographically close to Community Board 8. The process shows that tracts with many similarities to those in Community Board 8 are generally close to it- affluent census tracts are birds of a feather.

FIGURE 3 ABOUT HERE

The upper right corner of the SOM feature map represents the more affluent portions of the city (Figure 4). Selecting the neurons in the extreme upper right yields a geographic map that includes some of the more affluent parts of the city (including portions of the West Village, Chelsea, the Upper East and West sides, Forest Hills, Park Slope, Brooklyn Heights, and Riverdale). The portion of the SOM that is most distant from the upper right in the attribute space, the lower left, corresponds with tracts in northern Manhattan, Harlem, and the Bronx (Figure 5). The extreme lower left contains census tracts where over 50% of the population lives in poverty (as defined in the 2000 census). The tracts of the lower left do not group into clearly defined areas as well as those in the upper right. This suggests that these parts of the city that are most dissimilar to the affluent parts of the city exhibit less clustering than affluent areas, stated crudely, poor people are less clustered than rich people. In this example places with similar levels of attribute clustering show different levels of geographic clustering. The observation should be tempered by the U-matrix which shows that some differences between the upper right and lower left corners of the SOM.

FIGURE 4 ABOUT HERE

FIGURE 5 ABOUT HERE

Another approach to exploring and interpreting the trained SOM feature map is to select an area based on a single criterion, say census tracts where more than 90% of residents are not Caucasian, Figure 6 identifies census tracts that meet this particular criterion. Using the SOM feature map, we can find places that belong to the same classes as places where more than 90% of residents are not Caucasian (tracts mapped to the same neurons) thus highlighting areas that are similar. In the case of New York City, Figure 6 indicates that the latter areas are adjacent to zones of non-white concentration. The geographic pattern of places with a large non-white population is very similar to the pattern of these census tracts on the attribute map, there is a close link between the distribution of these tracts in geographic and attribute space. This approach paints a richer picture of the city, instead of using only a single criterion to identify similar places; we can now find areas of the city that are similar in many respects.

FIGURE 6 ABOUT HERE

Finally, one of the precepts of the human-ecologic approach that underlies geodemographics and urban factorial ecology is that populations sort themselves geographically to form socio-economically differentiated areal units or neighborhoods (Park & Burgess, 1925; Robson, 1969). This framework is important to the interpretation of census reporting districts. The degree to which census tracts fit this framework is an important question that while suited to the method reported here cannot be answered because the requisite disaggregate data is unavailable. The SOM method is a powerful tool to extract high-level structures of groupings of census tracts in the multidimensional attribute space. The comparison of patterns or structures in geographic space and attribute space is of interest as it sheds light onto the basic hypothesis of the human-ecologic approach to urban analysis. The relationship among and between map pairs can be assessed using a simple scaled measure of the average distance between observational units that are part of a subset of interest. This relative dispersion index can be formulated as

$$\frac{\sum_{i(case)} \sum_{j(case)} d_{ij} / N(case)}{\sum_{i(all)} \sum_{j(all)} d_{ij} / N(all)}$$

where d_{ij} is the Euclidean distance between observations i and j . To compare maps, we compute the average distance between all pairs of census tracts and all pairs of SOM buckets (neurons) that satisfy some criteria. We also compute the average distance between all pairs of census tracts and neurons that satisfy some criteria.

For example, following Figure 6 we selected all tracts where more than 90% of the population is self-identified as non-white. We then select the neurons of the self-organizing map that contain these tracts. We compute the average distance among all census tracts and among all neurons that satisfy this criterion. We then compute the relative dispersion by scaling these results by the average distance between all tracts and all neurons respectively. A relative dispersion greater than 1 indicates that the census tracts of interest are on average farther apart than all census tracts (more dispersed). Small numbers indicate that the tracts are on average closer together than all tracts (less dispersed). In the second column of Table 1 the small numbers for the community boards indicate that these are compact regions, as one would expect. The large number for tracts that are more than 90% white indicates that tracts which meet this criterion are widely dispersed throughout the city.

The third column on Table 1 reports the dispersion index for various criteria on the SOM feature map. The fourth column of Table 2 compares the two relative dispersion statistics; numbers near one indicate similar maps. High numbers in the fourth column indicate that the geographic distribution of observations is more dispersed than their attribute distribution. Low numbers in the comparison of map dispersion indicate less dispersion in attribute space relative to geographic space. All three socio-demographic criteria tested for dispersion point to dispersion in attribute space that very closely mirrors dispersion in geographic space these measures support widely observed patterns of racial segregation. Parts of the city that are more than 90% Caucasian are more spread out in physical and attribute space than places that are heavily African-American or minority. Supporting the discussion of Community Board 8 earlier in the paper one finds that its census tracts are highly clustered in attribute space. In general one finds that the very affluent parts of the city, for example, those tracts in the top 1% for income, are less diverse and more geographically concentrated than the lower income parts of the city.

Table 1. Relative Dispersions in Geographic Space and in Attribute Space.

Criteria	Relative Census Tract Dispersion	Relative Neuron Dispersion	Comparison of Relative Dispersions (Census Tract Dispersion / Neuron Dispersion)
90% Minority (Figure 6)	0.79	0.77	1.03
90% African-American	0.67	0.60	1.12

90% Caucasian	1.11	0.99	1.12
Manhattan Community Board 8	0.09	0.39	0.23
Brooklyn Community Board 6	0.12	0.68	0.17
Upper Right Corner of the SOM (Figure 4)	0.53	0.26	2.04
Lower Left Corner of the SOM (Figure 5)	0.84	0.26	3.23
Top Quartile for Median Household Income	1.14	1.03	1.12
Bottom Quartile for Median Household Income	0.87	0.96	0.91
Top Decile for Median Household Income	1.05	0.99	1.06
Bottom Decile for Median Household Income	0.87	0.99	0.88
Top 1% for Median Household Income	0.61	0.64	0.97
Bottom 1% for Median Household Income	0.91	0.87	1.05

Conclusions

Debates on the role of local demographic information in policymaking and governance are nearly as old as the United States. Representing the complexity of urban populations through cartography has been an area of inquiry since the 1890s. As data became more abundant and statistical techniques more refined, social area analysis and then factor analysis emerged. Modern geodemographic techniques have their roots in the analytic framework of the Chicago School and the methods pioneered by Eshref Shevky. Factor analysis and geodemographic techniques are limited in that they require researchers to describe the categories they have defined through labeling. This process of labeling has been critiqued in both factor analysis and geodemographics (Goss, 1995; Palm & Caruso, 1972). Self-Organizing maps are a novel approach to the problem of describing urban populations. They are novel in that when combined with geographic information systems they allow one to characterize the census tracts without the use of named categories. There have been relatively few applications of the Kohonen self-organizing map method to urban ecological analysis; we hope we have presented an accessible introduction to the utility of SOMs.

Self-Organizing maps share many of the limitations of factor analysis and geodemographics. When these techniques are applied to census divisions they must be interpreted with caution. Any analysis of census tracts in an urban area raises important questions about the nature of

tracts. Are tracts a meaningful unit of analysis? An analysis of census tracts is not an analysis of people and one must be careful to limit inference to scale of observation- any statements in this paper are about groups, not individuals. How important is the variability of populations within a tract to the overall classification scheme that results from a particular analytical approach?

The ability to visualize SOMs using commercial geographic information systems is limited. Interfaces between the SOM data mining method and GIS are not widely available however with lots of pointing and clicking or some simple scripting the connections can be made. The Geo Vista center at Penn State has developed tools for visualization of Self-Organizing maps in a geovisualization context (Takatsuka, 2001).

Finally, one of the most important aspects of using Self-Organizing maps in demographic analysis is variable selection. The absence of suitable theory to guide variable selection is a troubling reality; there is no current analogue to the Shevky-Bell hypothesis. Absent theoretical guidance the best a researcher can do is choose variables deemed important to the problem at hand. SOMs are an exploratory technique; they are not useful for confirming theory per se. Nor are they easily integrated into traditional statistical modeling techniques. While SOMs are subject to criticism because of their inability to extend urban theory, their ability to help people understand, the “particular circumstances of the community” make them a useful tool.

APPENDIX. Socio-Economic Variables

SQMILES	Area in square miles
POP100	Total population
HU100	Total housing units
POPDENS	Population density (POP100/SQMILES)
MALE_TOT	Total male population
FEM_TOT	Total female population
USCHLAGE	Population under school age, under 5 years
SCHLAGE	School age population, 5 to 17 years
MELDR_65	Elderly male population, 65 years and over
FELDR_65	Elderly female population, 65 years and over
ELDR_65	Elderly population, 65 years and over
PCT_USCA	Percent of population under school age, under 5 years
PCT_SCHA	Percent of school age population, 5 to 17 years
PCT_ELDR	Percent elderly population, 65 years and over
PCT_FEM	Percent female population
ENGLISH	English spoken at home, 5 years and over
SPANISH	Spanish spoken at home, 5 years and over
CHINESE	Chinese spoken at home, 5 years and over
RUSSIAN	Russian spoken at home, 5 years and over
ITALIAN	Italian spoken at home, 5 years and over
PCT_FORLAN	Percent foreign language spoken at home, 5 years and over
PCT_NATIVE	Percent native born
HU_OCC	Occupied housing units
HU_VAC	Vacant housing units
HU_OWN	Owner occupied housing units
HU_RENT	Renter occupied housing units
PCT_VACT	Percent vacant housing units
PCT_OWOC	Percent owner occupied housing units
MEDMOVED	Median year householder moved into housing unit
SAME1995	Population in same house in 1995
MEDRENT	Median contract rent quartile in dollars
PCTINCOME	Median gross rent as percent of household income in dollars
YRBUILT	Median year structure built
MEDVALUEOO	Median value for owner occupied housing units
HH_TOT	Total households reported
HH_POP	Total population in households
HH_AV_SZ	Average household size
HH_1PER	One person households
HH_FAM	Two or more person family households
HH_CH	Households with one or more people under 18 years
PCT_HHCH	Percent households with children
PCT_ALONE	Percent living alone
PCT_FAM	Percent in families
PCT_FHHH	Percent single female head of households
PCT_SMOM	Percent single mothers
TOT_FAM	Total families
POP_FAM	Total population in families

FAM_SIZE	Average family size
PCT_MARR	Percent of families married
PCT_MWC	Percent of families married with children
MED_HHI	Median household income
PERCAPITA	Per capita income in 1999
MEDINCOME	Total median earnings in 1999
PCT_POVERT	Percent below poverty level
PCT_UNEMP	Percent of workforce unemployed
PCT_DRIVE	Percent of workers driving to work
PCT_CAR	Percent of occupied housing units with vehicle available
PCT_PUB	Percent enrolled in public school (grades Pre-K to 12)
PCT_GRAD	Percent high school graduates, 25 years and over
WHITE	Number who self identify as only white (White alone)
BLACK	Black or African American alone
NATAMER	American Indian and Alaska Native alone
ASIAN	Asian alone
PACISLAND	Native Hawaiian and Other Pacific Islander alone
OTHER	Some other race alone
MULTIRACE	Two or more races
HISPANIC	Hispanic or Latino
NONHISP	Non-hispanic
ONE_NH	One race, non-hispanic
WHITE_NH	White alone, non-hispanic
BLACK_NH	Black or African American alone, non-hispanic
NATAM_NH	American Indian and Alaska Native alone, non-hispanic
ASIAN_NH	Asian alone, non-hispanic
PACIS_NH	Native Hawaiian and Other Pacific Islander alone, non-hispanic
OTHER_NH	Some other race alone, non-hispanic
MULTI_NH	Two or more races, non-hispanic
PCT_HISP	Percent Hispanic (HISPANIC/POP100)
PCT_WHITE	Percent White (WHITE/POP100)
PCT_MINOR	Percent Minority (1-WHITE/POP100)

Figure 1

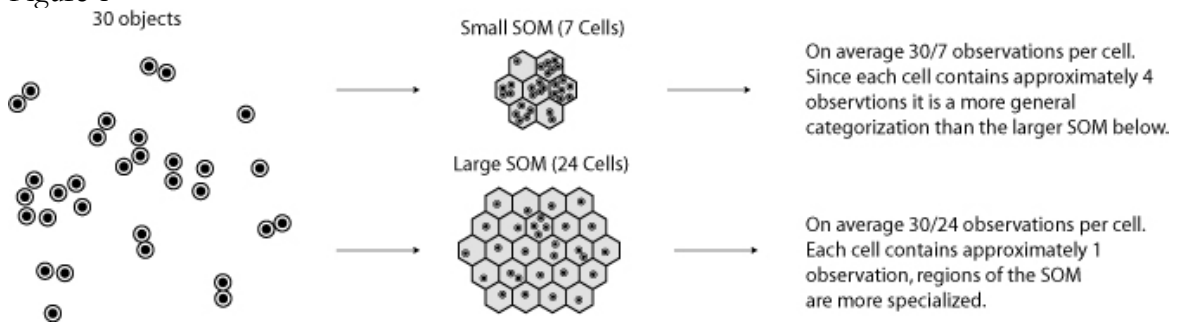


Figure 2

Unified Distance Matrix

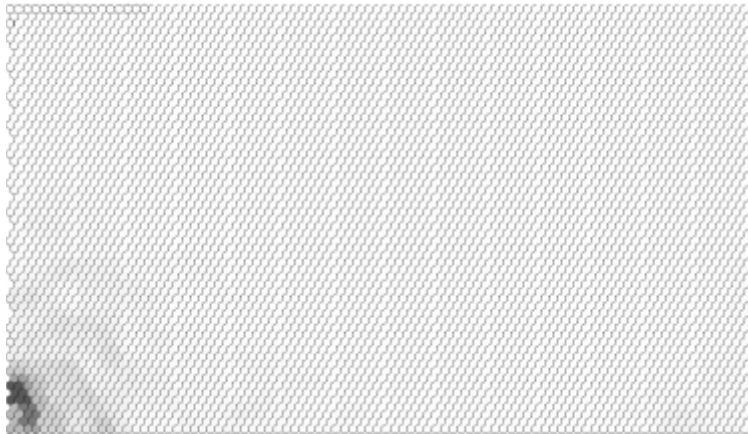
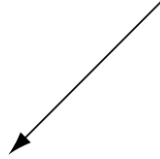
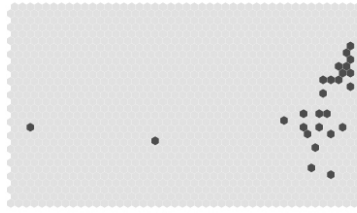


Figure 3

Community Board 8

Community Board 8 Feature Map



Tracts that are in the same class as Community Board 8

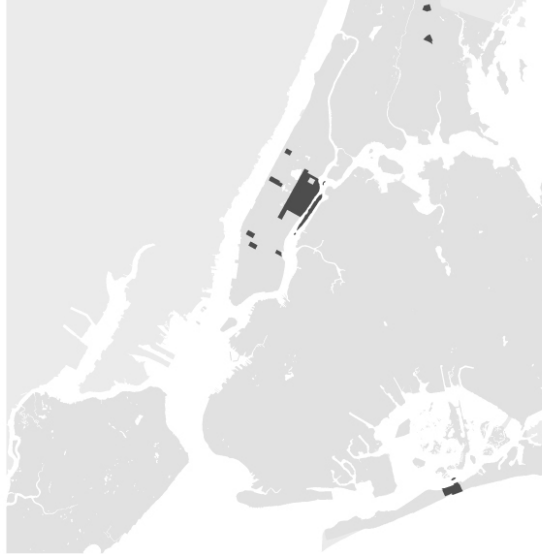
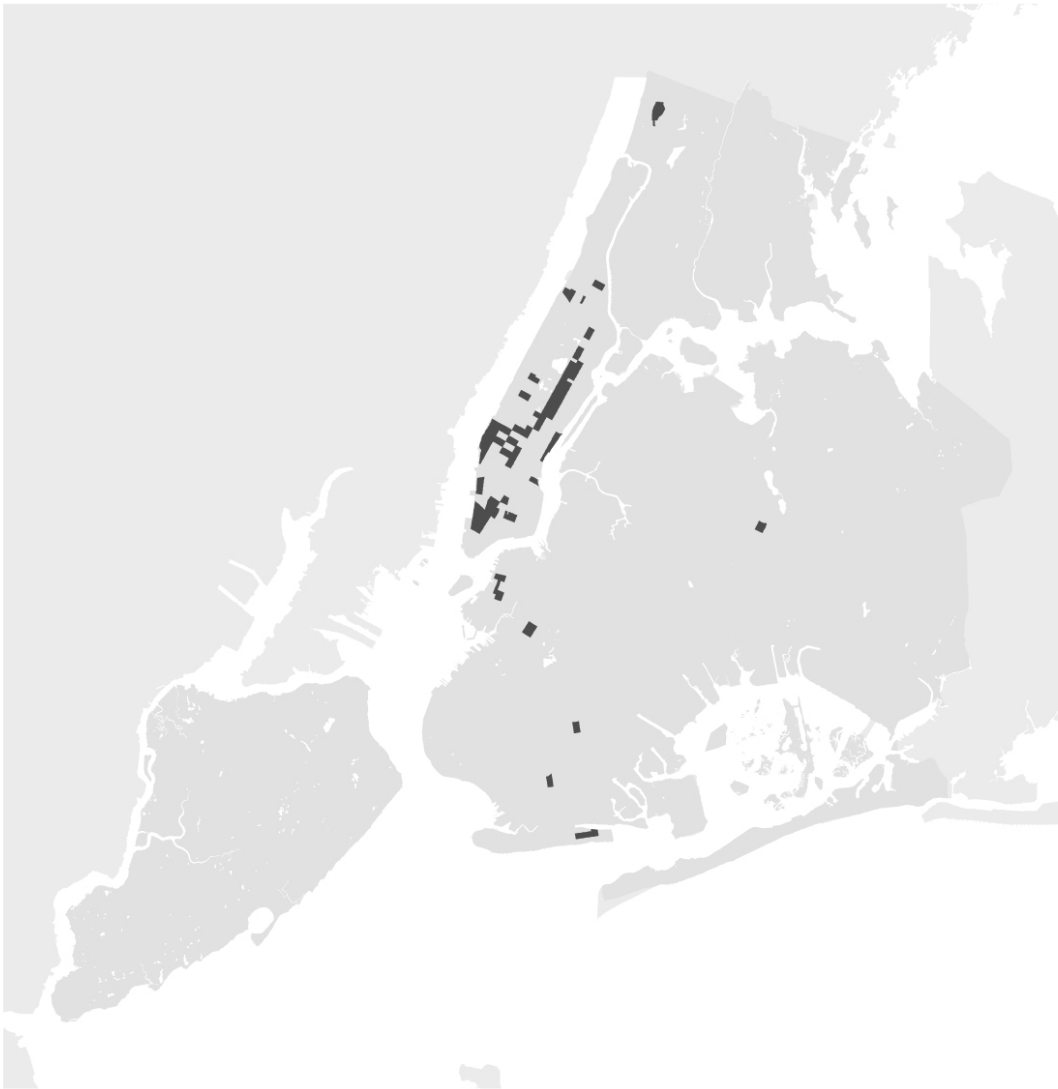


Figure 4



Census tracts in the upper right corner of the SOM

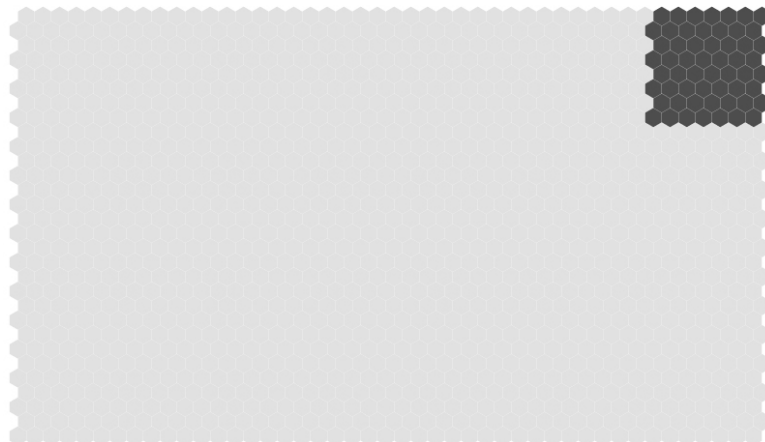
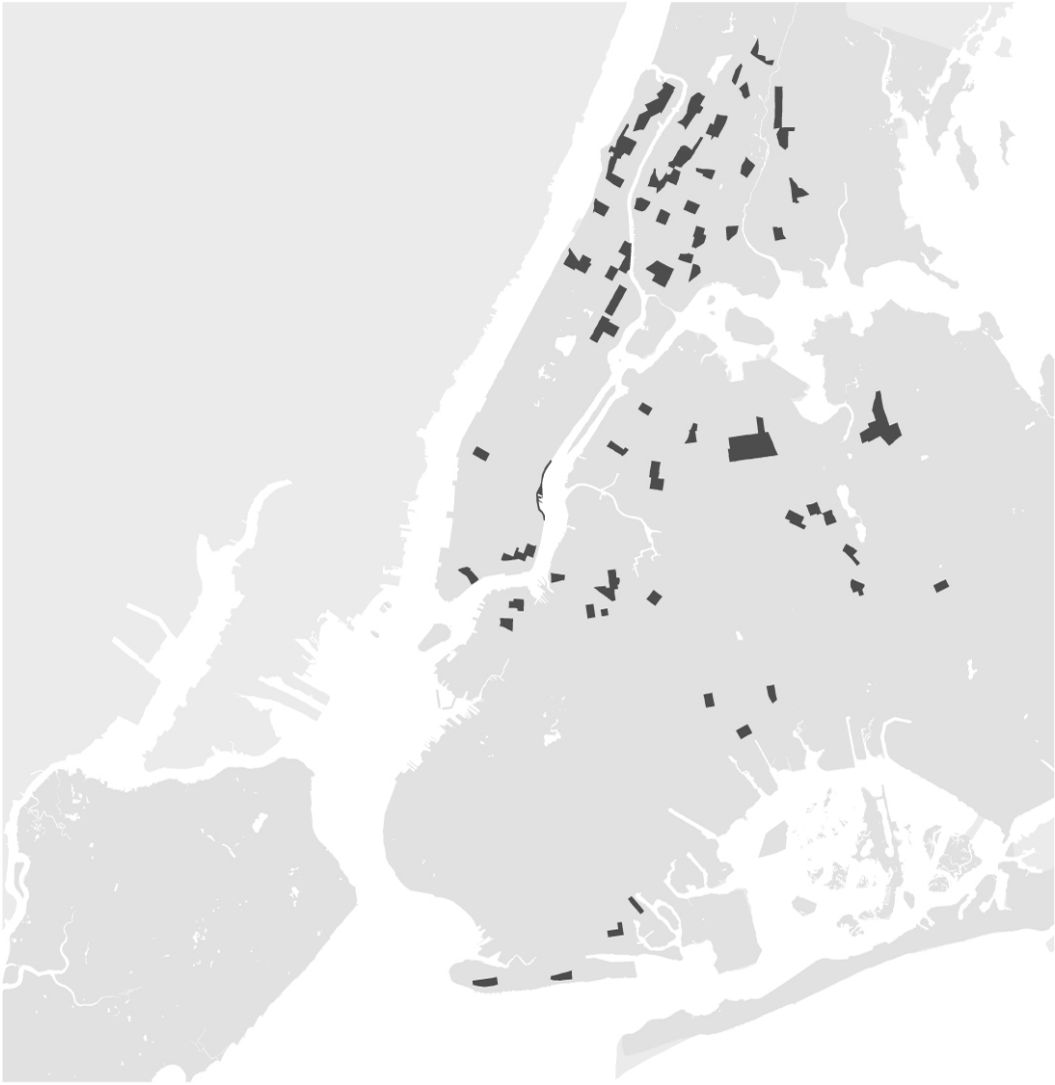


Figure 5



Census tracts in the lower left corner of the SOM

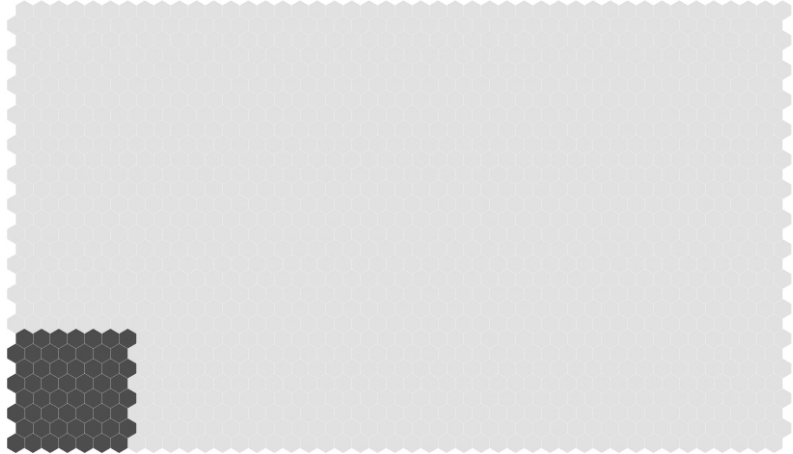
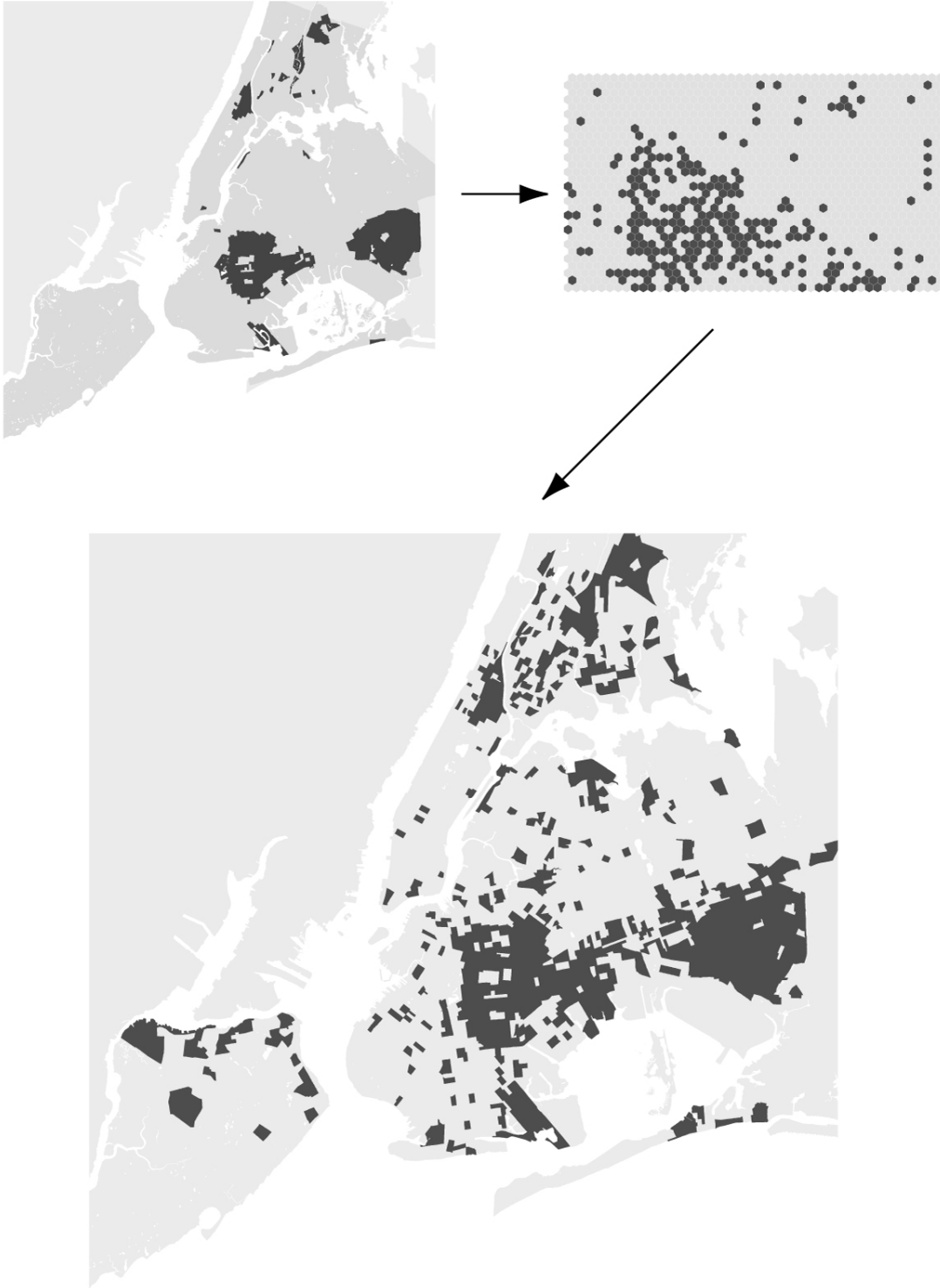


Figure 6

Census tracts where more than 90% of residents are not caucasian



Census tracts that map to the same regions of the feature map as those places that are more than 90% minority.

Bibliography

- Bacao, F., Lobo, V., & Painho, M. (2007). Applications of Different Self-organising Map Variants to Geographical Information Science Problems. In P. Agarwal, and A. Skupin (Eds.), *Self-organising Maps: Applications in Geographic Information Science*, Chichester: Wiley, forthcoming.
- Batey, P., & Brown, P. (1995). From Human Ecology to Customer Targeting: The Evolution of Geodemographics. In P. Longley and G. Clarke (Eds.), *GIS for Business and Service Planning*. Cambridge: GeoInformation International, pp. 77-103.
- Berry, B. J. L. (1971). Introduction: The Logic and Limitations of Comparative Factorial Ecology. *Economic Geography*, 47(2), 209-219.
- Booth, C. (1902). Map Descriptive of London Poverty, 1898-9: London School of Economics Charles Booth Online Archive.
- Cohen, P. C. (1981). Statistics and the State: Changing Social Thought and the Emergence of a Quantitative Mentality in America, 1790 to 1820. *The William and Mary Quarterly*, 38(1), 35-55.
- Curry, D. J. (1993). *The New Marketing Research Systems. How to Use Strategic Database Information for Better Marketing Decisions*. New York: Wiley.
- Goss, J. (1995). "We Know Who You Are and We Know Where You Live": The Instrumental Rationality of Geodemographic Systems. *Economic Geography*, 71(2), 171-198.
- Harris, R., Sleight, P., & Webber, R. (2005). *Geodemographics, GIS, and Neighborhood Targeting*. Chichester: John Wiley & Sons.
- Hatzichristos, T. (2004). Delineation of Demographic Regions with GIS and Computational Intelligence. *Environment and Planning B*, 31, 39-49.
- Himanen, V., Järvi-Nykänen, T., & Raitio, J. (1998). Daily Travelling Viewed by Self-Organizing Maps. In V. Himanen, P. Nijkmap, and A. Reggiani (Eds.), *Neural Networks in Transport Applications*. Ashgate. Aldershot, pp. 85-110.
- Holden, C. (2006). Hunter-Gatherers Grasp Geometry. *Science*, 311(5759), 317.
- Janson, C.-G. (1980). Factorial Social Ecology: An Attempt at Summary and Evaluation. *Annual Review of Sociology*, 6, 433-456.
- Johnston, R. J. (1976). Residential Area Characteristics: Research Methods for Identifying Urban Sub-areas - Social Area Analysis and Factorial Ecology. In D. T. Herbert and R. J. Johnston (Eds.), *Social Areas in Cities, v. I*, Chichester: Wiley, pp. 193-235.
- Kaski, S., & Kohonen, T. (1996). Exploratory Data Analysis by the Self-organizing Map: Structures of Welfare and Poverty in the World. In Refenes, Apostolos-Paul N. Refenes, Abu-Mostafa, Y., Moody, J., and Weigend, A. (Eds.). *Neural Networks in Financial Engineering*. Singapore: World Scientific, pp. 498-507.
- Kauko, T. (2005). Using the Self-Organising Map to Identify Regularities across Country-specific Housing-market Contexts. *Environment and Planning B*, 32 89-110.
- Kohonen, T. (1997). *Self-Organizing Maps*. Berlin: Springer.
- Kohonen, T., Hynninen, J., Kangas, J., & Laaksonen, J. (1995). The Self Organizing Map Program Package (Version 3.1). Espoo.
- Kurland, P., & Lerner, R. (Eds.). (1987). *The Founders Constitution*. Chicago: University of Chicago Press.
- Longley, P., & Clarke, G. (Eds.) (1995). *GIS for Business and Service Planning*. Cambridge: GeoInformation International.
- Miller, H., & Han, J. (Eds.) (2001). *Geographic Data Mining and Knowledge Discovery*. New York: Taylor & Francis.
- Openshaw, S. (1989). Neuroclassification of Spatial Data. In B. C. Hewitson, and R. G. Crane (Eds.), *Neural Nets: Applications in Geography*, Dordrecht: Kluwer Academic Publishers.
- Openshaw, S., & Wymer, C. (1995). Classifying and Regionalizing Census Tracts. In Openshaw, S. (Ed.). *Census User's Handbook*. Cambridge: GeoInformation International, pp. 239-270.
- Palm, R., & Caruso, D. (1972). Factor Labeling in Factorial Ecology. *Annals of the Association of American Geographers*, 62(1), 122-133.
- Park, R., & Burgess, E. (1925). *The City*. Chicago: University of Chicago Press.
- Rees, P. H. (1971). Factorial Ecology: An Extended Definition, Survey, and Critique of the Field. *Economic Geography*, 47(Supplement: Comparative Factorial Ecology), 220-233.

- Robson, B. T. (1969). *Urban Analysis: A Study of City Structure with Special Reference to Sunderland*. Cambridge: Cambridge University Press.
- Shevky, E., & Williams, M. (1972). *The Social Areas of Los Angeles: Analysis and Typology*. Westport: Greenwood Press.
- Skupin, A., & Agarwal, P. (2007). Introduction - What is a Self-Organizing Map. In P. Agarwal, and A. Skupin (Eds.), *Self-organising Maps: Applications in Geographic Information Science*, Chichester: Wiley, forthcoming.
- Skupin, A., & Fabrikant, S. I. (2003). Spatialization Methods: A Cartographic Research Agenda for Non-Geographic Information Visualization. *Cartography and Geographic Information Science*, 30(2), 95-119.
- Skupin, A., & Hagelman, R. (2003). Attribute Space Visualization of Demographic Change. In Proceedings of the 11th ACM International Symposium on Advances in Geographic Information Systems, New Orleans, LA, pp. 56-62.
- Skupin, A., & Hagelman, R. (2005). Visualizing Demographic Trajectories with Self Organizing Maps. *GeoInformatica*, 9(2), 159-179.
- Takatsuka, M. (2001). An Application of the Self-Organizing Map and Interactive 3-D Visualization to Geospatial Data. In Pullar, D. V., Proceedings of the 6th International Conference on GeoComputation, Brisbane, Australia, URL: <http://www.geocomputation.org/2001/papers/takatsuka.pdf>, accessed October 30, 2006.
- Thill, J.-C., Kretschmar, W., Casas, I., & X. Yao. (2007). " Detecting Geographic Associations in English Dialect Features in North America within a Visual Data Mining Environment Integrating Self-Organizing Maps. In P. Agarwal, and Skupin, A. (Eds.), *Self-organising Maps: Applications in Geographic Information Science*, Chichester: Wiley, forthcoming.
- Utley, G. (March 17, 2001). Cultural Diversity and America's High Schools. *CNN*.
- Villmann, T., Merenyi, E., & Hammer, B. (2003). Neural Maps in Remote Sensing Image Analysis. *Neural Networks*, 16(3-4), 389-403.
- Ward, D. (1969). The Internal Spatial Structure of Immigrant Residential Districts in the Late Nineteenth Century. *Geographical Analysis*, 1(4), 337-353.
- Ward, D. (1990). Social Reform, Social Surveys, and the Discovery of the Modern City. *Annals of the Association of American Geographers*, 80(4), 491-503.
- Webber, R. (1985). The Use of Census Derived Classifications in the Marketing of Consumer Products in the United Kingdom. *Journal of Economic and Social Measurement*, 13, 113-124.
- Webber, R. (2004). Designing Geodemographic Classifications to Meet Contemporary Business Needs. *Interactive Marketing*, 5(3), 219-237.
- Weiss, M. J. (1989). *The Clustering of America*. New York: Perennial Library.
- Witten, I., & Frank, E. (2005). *Data Mining: Practical Machine Learning Tools and Techniques, 2nd Edition*. San Francisco: Morgan Kaufmann.
- Yan, J., & Thill, J.-C. (2007). Visual Exploration of Spatial Interaction Data with Self-organizing Maps. In P. Agarwal, and A. Skupin (Eds.), *Self-organising Maps: Applications in Geographic Information Science*, Chichester: Wiley, forthcoming.