THE MODIFIABLE AREAL UNIT PROBLEM (MAUP) VIA CLUSTER ANALYSIS

METHODOLOGIES:

A LOOK AT SCALE, ZONING, AND INSTANCES OF FORECLOSURE IN LOS

ANGELES COUNTY


by


Matthew W. Davis


_____


A Thesis Presented to the
FACULTY OF THE USC GRADUATE SCHOOL
UNIVERSITY OF SOUTHERN CALIFORNIA
In Partial Fulfillment of the
Requirements for the Degree
MASTER OF SCIENCE
(GEOGRAPHIC INFORMATION SCIENCE & TECHNOLOGY)


May 2012

## Acknowledgements

# Table of Contents

# List of Tables

# List of Figures

**Abstract**

Spatial research has been plagued by the modifiable areal unit problem, or MAUP (Openshaw, 1979), for decades.  The MAUP can be broken down into two categories, the scale or aggregation effect and the zoning or grouping effect.  Recent advances in spatial science and technology have exacerbated the effects of the MAUP prevalent in many forms of research.  In this paper, data was obtained depicting instances of foreclosures in Los Angeles County (also the study extent) in 2006-2008.  This data was spatially joined to three sets of grids, or fishnets, covering Los Angeles County.  The data was also spatially joined to three additional datasets: the individual parcels of Los Angeles County and two aggregations of these parcels.  Five cluster analysis tools were used to analyze each of these seven total datasets.  Each permutation involved the five tools, seven datasets, and multiple pre-selected distance thresholds to test for scale and zoning effects of the MAUP. There were a total of 137 successful iterations illustrating the aforementioned permutations.  It was determined the MAUP is prevalent in this case study.  Suggestions are made in determining future actions to combat the effects of the MAUP.

**Chapter One – Introduction**

**1.1 Introduction and Rationale**

The advent of geographic information systems, or GIS, has propelled spatial

research and analysis to new heights.  Within the past few decades the acceleration of

computer technology coupled with the heightened prevalence of spatial data have forced

researchers to reexamine their assumptions regarding how best to address the inherent

flaws in such data.  These flaws include:  accuracy, scale, aggregation effects, and how

best to relate such data to other variables which may or may not be comprised of similar

spatial structures.

In order to understand the relationships between spatial datasets with different

extents it is necessary to aggregate data from one dataset into another.  One such method

is to aggregate point data into polygon data.  Methods of aggregation provide researchers

with the ability to cross-reference and analyze different types of data to see potential

spatial associations or relationships.  The boundaries with which these data are

aggregated tend to have arbitrarily defined extents.  For example, census tracts, city

boundaries, counties, states, and regularly gridded study areas are all artificially

constructed for various administrative or study-specific purposes.  These issues are a part

of the *scale effect* and *zoning effect* commonly found in multivariate spatial analysis.

These effects fall under the umbrella term known as the modifiable areal unit problem, or

MAUP.  The MAUP was coined by Openshaw & Taylor (1979) to describe the effects of

these two problems.  This thesis addresses MAUP issues and includes a discussion on

how the MAUP can be taken into consideration in order to prevent its effects from

biasing or leading to misinterpretation of results.  One of the ways this thesis proves the

existence of the MAUP is to utilize Local Indicators of Spatial Association, or LISA,

statistics (Dale et al., 2002; Ratcliffe and McCullagh, 1999; Lee and Rogerson, 2007;

Anselin, 1995).  This case study examines the relationship between foreclosure points in

Los Angeles County and how they relate to spatially joined aggregated polygons.  These

feature datasets share a similar extent though differ greatly among one another in terms of

polygon structure and placement.  Chapter Three includes more information on these

differences and what their long term implications are in spatial research.

This thesis is comprised of five chapters including this introductory chapter.

Chapter Two is a comprehensive Literature Review that examines research related to

spatial data queries and the conception of methodological frameworks, and critiques the

methods, results, and potential areas of improvement to the aforementioned research.

Chapter Three details the methodology used to examine the MAUP using LISA statistics

and includes a case study involving the relationship between foreclosures and units of

foreclosure aggregation as they relate to multiple cluster analysis methodologies.

Chapter Four details the results of the latter case study with maps, charts, and graphs of

the results of five different cluster methods.  Chapter Five interprets the LISA results

which lead into a discussion about how researchers need to include multiple permutations

to account for the effects of the MAUP.   Chapter Five also includes a broader

interpretation of the overall results of this thesis project.  A hypothetical example is also

included which illustrates how the effects of the MAUP need to be accounted for before

any actual research begins.

**1.2 Theoretical Framework and Questions**

Spatial research is limited to the dataset with the highest resolution or largest extent(s) set forth to be studied (Shriner, et al., 2006).  It is possible, though extremely difficult, to disaggregate such data into smaller scales or components for analysis with finer data (Gotway and Young, 2002).  This thesis will deal specifically with the issues that arise when point data is aggregated into larger areal units and the inherent concerns therein.

How thoroughly a given methodology has examined the effects of the MAUP indicates whether or not the researcher(s) have accounted for its effects.  Each variable in spatial research should require an exhaustive evaluation of the scale (or aggregation) and the zoning effects of the MAUP.  There is a trade-off in the level of uncertainty and error in the results/conclusions that a researcher must be willing to accept based on how thoroughly this analysis is performed.  These issues necessitate an awareness of how the MAUP issues bias and impact results based on the chosen methods of spatial calculation.  As completely as possible this thesis will answer the following questions related to foreclosure, aggregation, and zoning issues.  These questions are phrased as follows:

1) How does the modifiable areal unit problem (MAUP) affect the visual results of aggregate data in a multiscale environment via cluster analysis methods across multiple distance thresholds?

2) What is the best scale or aggregation unit used for running spatial analysis involving the use of local indicators of spatial association (LISA) statistics?

3) How do these effects impact the way in which a research methodology is constructed?

The impact of the MAUP can be a nuisance to any research that attempts to spatially link discordant variables across different spatial extents and scales (Getis and Franklin, 1987). One example of this situation includes instances of foreclosures, represented as point data, and crime rate statistics within city boundaries, represented as polygon data. One way to do a comparative measurement of the two variables would be to aggregate the point foreclosure data into the city boundaries to determine how foreclosure and crime rate correlate. However, such aggregation tends to generalize the foreclosure data. In this example, intra-city clusters of foreclosures and trends will unfortunately be eliminated through such aggregation methods. The fifth chapter of this thesis will discuss how best to approach this problem in one's research and provide insights as to how researchers might incorporate these approaches into their analysis.

**Chapter Two – Literature Review**

## 2.1 Problem Formulation

There has been moderate consideration given in the literature to the role zones and aggregation play in spatial datasets and subsequent analysis. Aggregation is defined herein as the inclusion of point/unit level data into larger arbitrarily defined extents that can be used for spatial analysis and statistical research (Pawitan and Steel, 2009). There are few studies which have attempted to account for or exonerate the effects of zonal data and spatial extent issues that plague many forms of spatial research. These sets of issues are known as the Modifiable Areal Unit Problem or MAUP (Openshaw and Taylor, 1979; Openshaw, 1984) and are broken up into two major effects, the scale effect (sometimes called the aggregation effect) and the zoning effect (Amrhein and Reynolds, 2006; Amrhein and Reynolds, 2007). This chapter details and critiques background literature which highlights the effect that aggregation has on analytical results. One of the goals of this chapter is to comprehensively articulate and summarize the necessary criteria for aggregate analysis.

In order to relate point-based feature data to other data in a spatial context it is often necessary to aggregate the attributes of those features into larger units in order to produce viable measures and statistics. In many cases the variable to compare is already aggregated into a larger extent, making point aggregation essential for analysis. The most common use of aggregates usually aligns with census blocks or tracts, zip codes, and city or county boundaries. Using established boundaries for aggregation provides the added benefit of linking one study's aggregates and conclusions to other studies with

similar extents and study areas. This literature review is intended to define and critique prior research while also providing the necessary background to synthesize a framework for the proposed methodology. In order to quantify the effect of the MAUP this research focuses on the different aspects of cluster analysis methodologies applicable to foreclosure point datasets based on total counts of foreclosure instances. The foreclosure data is utilized to analyze the effects of the MAUP by spatially relating point data to different feature boundaries and types.

This chapter is composed of four sections including this problem formulation section. Section 2.2 covers the history and background theories of aggregation, including the MAUP. Section 2.3 scrutinizes and critiques the methodologies and approaches employed by the articles summarized in Section 2.2. Finally, Section 2.4 reviews some pertinent research involving cluster analysis methodologies and compares them to each other. Chapter Five concludes this inquiry by focusing on possible methods for recognizing the impacts of the MAUP, and a discussion about future researchers might design their studies to lessen the effects of the MAUP.

## 2.2 Background Theories of the Components of the MAUP

The relevance of this study to the pertinent literature requires a comprehensive look at previous studies, in particular their contribution to understanding the MAUP, as well as the lens some used to compensate for the MAUP. Two primary problem paradigms are used to categorize and evaluate the MAUP which consists of the scale problem and the zoning problem. Jelinski and Wu (1996) define the scale problem as

"where the same set of areal data is aggregated into several sets of larger areal units, with each combination leading to different data values and inferences" and the zoning problem as "where a given set of areal units is recombined into zones that are of the same size but located differently, again resulting in variation in data values and…different conclusions".  In order to understand the reasons why certain methodological choices were made for this study it is necessary to evaluate the methodologies, conclusions, and schemas used by other researchers with similar goals.  By understanding preceding research and narrowing down extraneous avenues of thought this chapter is intended to provide clearer insight into the compositions of spatial research involving aggregate data and the MAUP.  Thus the following discussion provides the primary rationale for this study.

The term MAUP was coined by Openshaw and Taylor (1979).  A few studies were published in the 1980s and early 1990s (with the notable exception of Openshaw in 1984) which indicated that interest in the MAUP had tapered off (Amrhein, 1993).  The MAUP has recently become more prevalent and prominent in zone and scale research. The pronounced effect the MAUP has on research, particularly with different hotspot/cluster methodologies, warrants a level of scrutiny that will encourage future researchers to account for the bias, effects, and scope the MAUP.  This section notes several prominent studies (Table 2-1) and elaborates on the similarities and differences in results and conclusions between these studies as well as the corollary of the MAUP. Accounting for the MAUP provides the added benefit of eliminating an avenue of bias and thus potential criticism toward one's methodologies and conclusions.

Each of the studies discussed in this review have either attempted to measure the MAUP or compensate for it in some way in order to prevent zoning or scale/aggregation effects from warping statistical results and conclusions.  The following section details the studies focused solely on the scale effect; Section 2.4, the scale/aggregation effect; and lastly, section 2.5 covers research on scale and aggregation (Table2-1).

Table 2-1: Studies related to the MAUP described in this chapter (inclusive of some background literature).

| Author(s) | Scale or Aggregation Problem | Zoning Problem | Hotspot | Moving Window | Upscaling | OSA | OSU | COSP | Geocoded Data | LISA |
|---|---|---|---|---|---|---|---|---|---|---|
| Dark and Bram (2007) | X | | | | X | X | | | | |
| Ratcliffe and McCullagh (1999) | X | | X | X | | | | | X | X |
| Shriner, Wilson, and Flather (2006) | X | | X | | | | | | | |
| Mu and Wang (2008) | X | | | | | | | | X | |
| MacEachren (1982) | | X | | | | | | | | |
| Bhati (2005) | | X | | | | | | | X | |
| Hipp (2007) | | X | | | | | | | | |
| Hayward and Parent (2009) | X | X | | | | | | | | |
| Jelinski and Wu (1996) | X | X | | X | X | | | | | |
| Gotway and Young (2002) | X | X | | | X | | | X | | |
| Nakaya (2000) | X | X | | | | | | | | |
| Hay, Marceau, Dube, and Bouchard (2001) | X | X | | | X | X | X | | | |
| Pawitan and Steel (2008) | X | X | | | | | | | | |
| Tagashira and Okabe (2002) | X | X | | | | | | | | |
| Chainey, Thompson, And Uhlig (2008) | X | X | X | | | | | | X | |
| Gatrell, Bailey, Diggle, and Rowlingson (1995) | | | | X | | | | | | |
| Tagashira and Okabe (2002) | | | X | | | | | | | |
| Chainey, Thompson, And Uhlig (2008) | | | | X | | | | | | |
| Gatrell, Bailey, Diggle, and Rowlingson (1995) | | | | | | | | | | |
| Fotheringham (1989) | X | | | | | | | | | |
| Fotheringham and Wong (1991) | X | X | X | | | | | | | |
| Lentz, Blackburn, and Curtis (2011) | X | X | | | | | | | | |
| Amrhein (1993) | X | X | | | | | | | | |
| Amrhein and Reynolds (1996) | X | X | | | | | | | | |
| Amrhein and Reynolds (1997) | X | | | | | | | | | |
| Rushton and Lolonis (1996) | | | | X | | | | | | |
| Rushton (1998) | | | | X | | | | | | |

**2.3 The Scale Effect**

The scale problem, as defined previously (Jelinski and Wu, 1996), concerns how the scale of any given data impacts the way in which they are analyzed and interpreted. When moving from a fine scale geographic extent to one that excludes variation, certain types of error inevitably follow. Dark and Bram (2007), Ratcliffe and McCullagh (1999), and Shriner, et al. (2006) each attempt to approach the scale problem with different methodologies. Scale is largely chosen by the extents and scale of the input data for any evaluation. In many cases this involves the structure of a raster and the 'cells' that constitute the data composites to be analyzed. Dark and Bram (2007) emphasize the scale effect as seen in physical geography. Ratcliffe and McCullagh (1999) use hotspots of crime occurrences to illustrate the scale effect. Shriner, et al (2007) use five grain sizes, or raster extents, on richness hotspots for conservation planning regarding species reserve network. Mu & Wang (2008) analyze different ways of mitigating the scale effect.

Dark and Bram (2007) noted that the natural sciences were quick to recognize the natural scales of many geographic and ecological phenomena. Determining the scale with which to study a particular phenomenon or process requires a balance between the scales of the different available data and considerations for future cross analyses which utilize similar scales and extents. Dark and Bram's (2007) study focused on creating flow line models based on different Digital Elevation Models, or DEMs, of the same area. The coarsest, or lowest scale resolution data, produced rough flow lines based on 30 meter DEMs. When overlain on finer scale DEMs (10 meter and 1 meter,

10

respectively) one may visually asses how the scale problem can warp resultant data when assessing the topographical features of the flow lines. They also discuss possible solutions to the scale effect of the MAUP that are highly relevant to the methodology and results of this thesis.

Ratcliffe and McCullagh (1999), on the other hand, focus their study on hotspot analysis and LISA statistics regarding instances of crime. They compare actual crime statistics to the perceived zones of concentrated criminal activity as reported by police officers across three police force subdivisions. They account for and mitigate the scale effect of the MAUP by using a 'moving windows' approach to analyze the data. This particular approach is comprised of a "technique [that] applies a movable sub-region (usually a circle) over the entire study area to measure dependence in subsets of the study area" (Ratcliffe and McCullagh, 1999). The 'moving window' approach allows for the data to be analyzed without aggregation or the scale of the data creating a conflicting framework for analysis (Rushton and Lolonis, 1996; Rushton, 1998, Ratcliffe and McCullagh, 1999). The results of Ratcliffe and McCullagh's (1999) study illustrate the way in which hotspot analysis and the moving window technique can help law enforcement agencies concentrate their resources on endemic areas proven to be statistically significant hotbeds of crime. The moving window approach was pioneered by Rushton and Lolonis (1996) and Rushton (1998).

Shriner, et al (2006) performed their analysis using five different raster extents on species richness hotspot data to determine the spatial overlap of generated reserve networks. Richness hotspot reserves are comprised of areas where species diversity is

11

delineated by hotspot concentrations of specific species. Through multiple iterations of calculations they determined that "the number of species represented in richness hotspot reserves increased as grid cell size decreased" (Shriner, et al., 2006). This particular study clearly illustrated the scale effect of the MAUP since reserve network configurations varied greatly across different grain sizes, which were meant to highlight how the initial scale of input data affects the possible outcomes across any spatial analysis. Shriner, et al. (2006) also noted that thoughtful consideration is required when conceiving how fine-scale data may engender fragmentation and coarse-scale data may lose essential accuracy.

Unlike the previously mentioned research, Mu and Wang (2008) attempt to formulate a practical program for implementing a modified scale-space clustering method, called MSSC. This program is intended to "account for both attribute homogeneity and spatial contiguity" (Mu and Wang, 2008). To test this new technology they considered a case study involving homicide rates in Chicago in 1990. Using the MSSC method they develop and aggregate various tracts of land within the city of Chicago into groups containing fewer and fewer units as the calculations progress (continuous aggregation). Mu and Wang (2008) "name the process of repeatedly running a clustering method toward the formation of one single cluster a *converging effect*." Their results clearly indicate the presence of the MAUP, especially when applying this methodology to many different zones. This indicates how each process of a continuous aggregation yields a dataset that is increasingly generalized and homogenous when compared to its predecessor(s).

**2.4 The Zoning Effect**

The other category of the MAUP involves the zoning effect. This effect is marked by the assumption that information aggregated to a specified area is true for every subset of that area and that the resultant statistics of any calculations will change if the zones change although the total areal extent remains constant. This effect causes loss of the precision and accuracy of point-data. For example, this loss of precision would then favor area-based data representation which makes for easier interpretation and calculation (Williamson, et al., 2001). There are three additional articles in the following discussion that exemplify the zoning effect. These articles include analysis in choropleth map accuracy (MacEachren, 1982), different levels of neighborhood aggregate structure analysis and crime (Hipp, 2007), and using local levels to determine crime rates (Bhati, 2005).

At the time of publication of MacEachren's (1982) research, GIS technology did not exist to readily critique different types of aggregate structures in order to determine the areal units best suited for any given research. Despite this fact MacEachren (1982) constructed a methodology around individual noncontiguous enumeration units representing different counties across the United States. These units were overlaid on a regular grid of points to measure compactness, data distribution variability, and unit size. MacEachren (1982) concluded that unit size and enumeration are significant factors that need to be considered for any choropleth map creation, since subsequent interpretation is heavily dependent upon these factors. It can be inferred from MacEachren's (1982)

research that scale and zonal measures impact not just data structures in analysis and representation but also how information is visually assessed and interpreted.

While MacEachren (1982) took a broader view on unit configuration, Hipp (2007) proffers a detailed discussion on the common problems associated with selecting what constitutes 'neighborhoods' and how aggregating information with different zonal locations may emphasize or obscure natural clusters of information. Hipp (2007) then looks at various demographic characteristics as they relate to different criminological paradigms to see how aggregations with different zones impact results and conclusions drawn from them. Hipp (2007) concludes that there is no universally accepted "appropriate" zone of aggregation that all researchers must follow. Rather a researcher must take all possible zonal levels and scales into careful consideration in order to find the one that is most appropriate for a particular study. Hipp (2007) also presents the argument that reconciling how different people subjectively view the boundaries of a given area can result in vast differences in how neighborhood boundaries are constructed and analyzed. For example, how a city delineates a neighborhood boundary and how a citizen living in a neighborhood delineates that same neighborhood boundary can and do result in different definitions for areas and thus issues related to areal definitions.

Lastly, Bhati (2005) focuses his research on determining an analytical framework to approach rare crimes via areal aggregation at different local levels. Data from the census tract level was aggregated to neighborhood clusters for the analysis. Much like Hipp (2007), Bhati (2005) utilizes various demographic variables in his study. In this study the MAUP is only touched upon briefly, though the aggregate problem is shown to

be prominent in the resultant figures. Bhati (2005) goes on to assess that the two aggregate boundaries used in his analysis (census tract and neighborhoods) yielded statistically significant results across both sets with certain variables, while other variables had a statistical significance that fluctuated across both aggregation levels.

## 2.5 The MAUP (Both Effects)

Few studies look at both aspects of the MAUP, preferring instead to focus either on the scale effect or zoning effect. Seven primary articles were reviewed that deal with both facets of the MAUP. These studies deal with a broad range of scale problems and zoning problems, from poverty in Pennsylvania (Hayward and Parent, 2009) to upscaling landscape analysis (Hay et al, 2001), to incidence rate maps (Nakaya, 2000). The approaches that each researcher takes yield varying conclusions across a broad spectrum of research. The following paragraphs summarize the highlights of each article and provide a brief overview of the major conclusions.

Jelinski and Wu (1996) look at the MAUP as it relates to landscape ecology. Using Normalized Difference Vegetation Index, or NDVI, they were able to aggregate the units to test how the MAUP impacts spatial autocorrelation. Jelinski and Wu (1996) created an artificial example of the MAUP effects which was reproduced by Dark and Bram (2007) and seen here in Figure 2-1. In their analysis of landscape data Jelinski and Wu (1996) aggregated data from the finest scale (1x1), basic units, to coarsest scale (225x225). To measure the zoning effect they then used two aggregation procedures at different scales but "with an equal number of pixels per zone" (Jelinski and Wu, 1996).

To test for spatial autocorrelation they used Moran's I statistic and Geary's C statistic.

The scale effect of the MAUP can be seen when comparing the results of different spatial

autocorrelation operations which utilize different scales (Fotheringham, 1989;

Fotheringham and Wong, 1991). Jelinski and Wu (1996) go on to conclude that "from a

hierarchical stand point of view, the MAUP is not really a 'problem'…it may reflect the

'nature' of the real systems that are hierarchically structured" (Jelinski and Wu, 1996).

The determination that hierarchical structure is important in analyzing the MAUP and in

mitigating its effects will be discussed in more detail in section 2.6.

**Effects of Aggregation**

| a | | | | | b | | | c | |
|---|---|---|---|---|---|---|---|---|---|
| 2 | 4 | 6 | 1 | | 3 | 3.5 | | 3.75 | 3.75 |
| 3 | 6 | 3 | 5 | | 4.5 | 4 | | | |
| 1 | 5 | 4 | 2 | | 3 | 3 | | 3.75 | 3.75 |
| 5 | 4 | 5 | 4 | | 4.5 | 4.5 | | | |

| mean= 3.75 | mean= 3.75 | mean= 3.75 |
|---|---|---|
| variance= 2.60 | variance= 0.50 | variance= 0.00 |

**Effects of Zoning Systems**

| d | | | | | e | | | | f | |
|---|---|---|---|---|---|---|---|---|---|---|
| 2.5 | 5 | 4.5 | 3 | | | | | | 4 | 1 |
| | | | | | 2.75 | 4.75 | 4.5 | 3 | | |
| 3 | 4.5 | 4.5 | 3 | | | | | | 4 | 3.7 |

| mean= 3.75 | mean= 3.75 | mean= 3.17 |
|---|---|---|
| variance= 0.93 | variance= 1.04 | variance= 2.11 |

Figure 2-1: Artificial example of the aggregate effect and the scale effect on average and variance (Jelinski and Wu, 1996)

Hayward and Parent (2009) concentrate on quantifying the effects of the MAUP on poverty rates in Pennsylvania.  To learn about the scale effect Hayward and Parent (2009) created choropleth maps at different scales and zonal configurations which resulted in the use of a Kolmogorov-Smirnov (K-S) test.  To investigate the zoning effect they used varying district designs to compare rates between different aggregated areas.  Aside from the statistical proofs of the MAUP they also generated several choropleth maps of the results that showed different poverty concentrations that varied greatly in location and distribution.  They were careful to note how these proofs validated not only the existence of the MAUP but how it affects the public interpretation of poverty locations.

Hay et al. (2001) described their study and how they hoped to use hierarchically nested data structures to analyze dominant landscape features via the Object-Specific Analysis, or OSA.  They worked on satellite imagery of a small 36 ha2 area of an island in Canada.  This imagery was scaled (by increasing grain size, or cell resolution) to different scales and zones (resolutions) and calculated against the OSA and the Object-Specific Upscaling, or OSU, techniques to determine which was more apt at negating or controlling for the MAUP.  Raster grids are composed of regularly gridded cells, each with its own value or set of values.  The OSU is "unique, in that it incorporates an explicit multi-resolution (i.e., hierarchical) sampling and evaluation" of pixels within the larger context of its higher-grained extent (Hay et al., 2001).  Multiple upscaling iterations were used on the base imagery which led them to conclude that, of the two techniques, OSU is the one that is most appropriate for this type of study.  They go on to

suggest that "multiscale image-object thresholds are often far more 'fuzzy' or less discrete than the term threshold commonly implies" (Hay et al. 2001).

Nakaya (2000) studied the geographic relationship between a standard mortality ratio (SMR) and a designation of 262 basic areal units (BSUs) for Tokyo City, Japan. Incidence rate maps were generated based on several socioeconomic ratios comprising Tokyo. A comprehensive analysis was then performed utilizing variables and the geographic extent of the city and its comprising BSUs. Nakaya (2000) discusses possible factors that may skew the results concerning population density distributions, and those impacts on the final statistics. Nakaya (2000) concludes that "it is possible to find the proper resolution of maps depending on the size of the data" to adequately understand SMR rates and distributions.

Pawitan and Steel (2009) attempt to explore the MAUP using various statistical simulations on Australian Census data. The results of this analysis indicate strong evidence of MAUP, despite utilizing a single unaltered zone scheme. This thesis project will utilize a similar analysis schema to Pawitan and Steel (2009). Pawitan and Steel (2009) used two dataset aggregations while this thesis project will utilize four dataset aggregations. Chapter Three of this project discusses the procedures involved in the creation of the datasets used in this thesis.

Tagashira and Okabe (2002) attempt to study the MAUP from the specific framework of fixed predetermined point analysis. Their regression model is based on studying income and household location as a function of distance from the predetermined point of a business district. This analysis also attempts to deal with the ecological fallacy

problem in which "results calculated from aggregate level data are unlikely to correspond to those from individual level data" (Tagashira and Okabe, 2002). Through the use of efficiency analysis of slope coefficients it was determined that the MAUP was prevalent in this research area, and that the range of potential analyses roughly determines how best to measure the ecological fallacy problem and the MAUP.

Chainey et al. (2008) focused on spatial analysis of crime instances and hotspot analyses of this variable. The primary critique Chainey, et al. (2008) performs was an analysis of five different hotspot analysis techniques to determine which is more accurate in terms of estimating where actual future crimes will occur. They studied the London, England area with two additional adjacent areas that fell within the jurisdiction of the London Metropolitan Police Force. The results indicated that kernel density estimation, or KDE, was consistently the most accurate technique available for use with hotspots. Visualizations of the different hotspot techniques lend credence to their assertions that hotspot analyses have not been as rigorously tested in past studies as they should be. They also created a new prediction accuracy index, or PAI, as part of their validity testing of hotspot analysis. It is important to note at this point that this thesis project attempts to study hotspot techniques, though only under the paradigm of MAUP analysis and confirmation, rather than determining the primacy or superiority of any one hotspot technique.

Unlike the previously mentioned studies, Gotway and Young (2002) focus their research on how best to approach the MAUP through statistical methods and evaluation of best-fit approaches based on the scale of the data and possible aggregations therein.

Gotway and Young (2002) present a strong case for how the MAUP is merely a facet of the larger Change of Support Problem, or COSP, that permeates spatial research. They go on to define 'support' as "the size or volume associated with each data value, but the complete specification of this term also includes the geometrical size, shape, and spatial orientation of the regions associated with the measurements" (Gotway and Young, 2002). The aggregation of any data into different scales or areal units, in turn, creates entirely new variables that are only distantly related to the properties of the original variables. Rather than apply their methods to a particular case study, they suggest various statistical methods for dealing with particular COSPs.

## 2.6 Critiques

The following discussion deconstructs the issues involved with analyzing the MAUP as it relates to the different methodologies inherent in each study. Gotway and Young (2002) present examples of observations of the COSP and the nature of its processes. The four types of COSPs include point, line, area, and surface. The nature of each of the four COSPs may be quite different and directly affects the types of possible analyses that can be performed. Thus COSP can be considered the umbrella set of problems that the ecological fallacy problem and the MAUP are part of.

Of the literature reviewed some studies analyzed only one aspect of the MAUP and failed to take into consideration the effects of the other aspects of the MAUP within each analysis. Dark and Bram (2007), for example, make use of flow line analysis based on fine and coarse grain resolution imagery. The aggregation effect does not play a role

in this type of research unless other factors are being measured against that imagery or if the scale of the imagery changes across multiple dimensions. Dark and Bram (2007) do not take the scale effect into account in their examination.

Gotway and Young (2002) also noted that hierarchical structures are important for disaggregating clusters of data into more manageable units. This study reasonably asserts that such hierarchical disaggregation may make false assumptions about the accuracy of the areal data being dispersed. "These assumptions are much more than mathematical and computational assumptions…they can (often surreptitiously) result in unusual or unrealistic cross-scale relationships" (Gotway and Young, 2002). Shriner, et al. (2006) made a similar note with the comment that "many authors that use coarse-grained data either implicitly assume patterns detected using coarse grained data reflect those at other grains (Rahbek and Graves 2000, Larsen and Rahbek 2004) or fail to discuss potential consequences of extrapolating results to other scales" (Shriner, et al. 2006). More work needs to be done to ask these kinds of questions in order to critically focus on the framework of potential analyses that may deal with the MAUP and its impacts on analytical methods and conclusions.

Jelinski and Wu (1996) and Ratcliffe and McCullagh (1999) both primarily deal with the 'moving windows' approach in their analyses to the MAUP. The moving window technique is defined as applying "a movable sub-region (usually a circle) over the entire study area to measure dependence in subsets of the study area and is particularly suited to crime hotspot detection" (Ratcliffe and McCullagh, 1999). Utilizing the moving windows approach is similar to the way in which this thesis project

is constructed since it uses a distance threshold on each areal unit in some of the cluster

analysis methods. More on this methodology can be found in Chapter Three.

**Chapter Three – Methodology**

There are many potential vehicles of analysis available to those attempting to study the MAUP in spatial research. This thesis uses foreclosure events in Los Angeles County as the primary spatial variable. A total of 55,631 geographic locations pinpoint instances of foreclosure over a three year period (2006-2008). The use of these points, as well as artificially constructed boundaries, sheds a great deal of light on how pervasive the MAUP is in spatial research and its accompanying conclusions. This chapter is divided into five main sections: data gathering and acquisition, data analysis procedures, data preparation, cluster analysis methods, and justification. Each section clarifies the use of specific cluster analysis tools, data preparation techniques, distance thresholds, and geographically aggregated boundaries (Getis and Ord, 1992).

### 3.1 Data Gathering and Acquisition

The Los Angeles County Assessor's Office provided the original version of the point dataset used to depict the instances of foreclosures in LA County for the given time frame and the parcel boundary. All subsequent files used in this study were generated independently as a result of spatial tool processes which combined, aggregated, merged, joined, dissolved, or calculated cluster statistics on dataset features.

### 3.2 Data Analysis Procedure

The procedures involved in this research uses three main dimensions of analysis to understand the effects of the MAUP. There are five different cluster analysis tools

which are used on seven different datasets (Mueller-Warrant, Whittaker and Young III, 2008). Each of these datasets is then broken down into another subset via the distance threshold variable (Table 3-1). This variable indicates how far away a given point or polygon will 'look' when performing statistical analyses to determine cluster properties. Each tool, with the exception of Ripley's K, utilizes this threshold variable. Multi-Distance Spatial Cluster Analysis (Ripley's K-function) tool uses regularly spaced distance bands as a property of the tool itself. The procedures and model runs for this project used eight different distance thresholds: 164 meters, 300 meters, 500 meters, 750 meters, 1 kilometer, 2 kilometers, 5 kilometers and 10 kilometers. The first distance threshold, 164 meters, was determined to be the lowest acceptable distance threshold for a hotspot analysis to run on the original foreclosure point dataset. Since the foreclosure points are the base level data for each subsequent aggregation and tool iteration, this minimum and subsequent thresholds were maintained across the other geographically aggregated features. It is important to note that many processes were unable to utilize a distance threshold of 164 meters or more due to limitations of analysis software (Esri ArcGIS). Table 3-1 provides a complete reference as to which feature classes could run each cluster analysis tool and at which distance thresholds they would function. This table functions as a reference to the results maps in Chapter Four as well as the maps in the Appendices.

| Analysis Method | Iteration Results | |
| :---: | :---: | :---: |
| | **Foreclosure Points** | |
| | *Distance* | *Yielded Results* |
| **High/Low Clustering (Getis-Ord General G) HTML** | 164m | Yes |
| | 300m | Yes |
| | 500m | Yes |
| | 750m | Yes |
| | 1km | Yes |
| | 2km | Yes |
| | 5km | Yes |
| | 10km | Yes |
| **Multi-Distance Spatial Cluster Analysis (Ripley's K Function)** | N/A | Yes |
| **Spatial Autocorrelation (Moran's I) HTML** | 164m | N/A |
| | 300m | N/A |
| | 500m | N/A |
| | 750m | N/A |
| | 1km | N/A |
| | 2km | N/A |
| | 5km | N/A |
| | 10km | N/A |
| **Cluster and Outlier Analysis (Anselin Local Moran's I)** | 164m | N/A |
| | 300m | N/A |
| | 500m | N/A |
| | 750m | N/A |
| | 1km | N/A |
| | 2km | N/A |
| | 5km | N/A |
| | 10km | N/A |
| **Hotspot Analysis (Getis-Ord Gi*)** | 164m | N/A |
| | 300m | N/A |
| | 500m | N/A |
| | 750m | N/A |
| | 1km | N/A |
| | 2km | N/A |
| | 5km | N/A |
| | 10km | N/A |

Table 3-1:  Guide to model iterations and results.

*Indicates the distance threshold was too low for the model iteration to run.

**Indicates the available hardware was unable to compute the results.

| Analysis Method | Iteration Results | | |
| --- | --- | --- | --- |
| | | Parcels | Parcels2 |
| | *Distance* | *Yielded Results* | *Yielded Results* |
| **High/Low Clustering (Getis-Ord General G) HTML** | 164m | * | * |
| | 300m | * | * |
| | 500m | * | * |
| | 750m | Yes | Yes |
| | 1km | Yes | Yes |
| | 2km | ** | Yes |
| | 5km | ** | Yes |
| | 10km | ** | Yes |
| **Multi-Distance Spatial Cluster Analysis (Ripley's K-Function)** | N/A | ** | Yes |
| **Spatial Autocorrelation (Moran's I) HTML** | 164m | * | * |
| | 300m | * | * |
| | 500m | * | * |
| | 750m | Yes | Yes |
| | 1km | Yes | Yes |
| | 2km | Yes | Yes |
| | 5km | Yes | Yes |
| | 10km | Yes | Yes |
| **Cluster and Outlier Analysis (Anselin Local Moran's I)** | 164m | * | * |
| | 300m | * | * |
| | 500m | * | * |
| | 750m | Yes | Yes |
| | 1km | Yes | Yes |
| | 2km | Yes | Yes |
| | 5km | Yes | Yes |
| | 10km | ** | Yes |
| **Hotspot Analysis (Getis-Ord Gi*)** | 164m | * | * |
| | 300m | * | * |
| | 500m | * | * |
| | 750m | Yes | Yes |
| | 1km | Yes | Yes |
| | 2km | Yes | Yes |
| | 5km | Yes | Yes |
| | 10km | Yes | Yes |

Table 3-1, continued:  Guide to model iterations and results.

| Analysis Method | Iteration Results | | |
|---|---|---|---|
| | | Parcels3 | Fishnet1 |
| | *Distance* | *Yielded Results* | *Yielded Results* |
| **High/Low Clustering (Getis-Ord General G) HTML** | 164m | * | * |
| | 300m | * | Yes |
| | 500m | * | Yes |
| | 750m | Yes | Yes |
| | 1km | Yes | Yes |
| | 2km | Yes | Yes |
| | 5km | Yes | Yes |
| | 10km | Yes | Yes |
| **Multi-Distance Spatial Cluster Analysis (Ripley's K-Function)** | N/A | Yes | ** |
| **Spatial Autocorrelation (Moran's I) HTML** | 164m | * | * |
| | 300m | * | * |
| | 500m | * | * |
| | 750m | Yes | Yes |
| | 1km | Yes | Yes |
| | 2km | Yes | Yes |
| | 5km | Yes | Yes |
| | 10km | Yes | Yes |
| **Cluster and Outlier Analysis (Anselin Local Moran's I)** | 164m | * | * |
| | 300m | * | * |
| | 500m | * | * |
| | 750m | Yes | Yes |
| | 1km | Yes | Yes |
| | 2km | Yes | Yes |
| | 5km | Yes | Yes |
| | 10km | Yes | Yes |
| **Hotspot Analysis (Getis-Ord Gi*)** | 164m | * | * |
| | 300m | * | Yes |
| | 500m | * | Yes |
| | 750m | Yes | Yes |
| | 1km | Yes | Yes |
| | 2km | Yes | Yes |
| | 5km | Yes | Yes |
| | 10km | Yes | Yes |

Table 3-1, continued:  Guide to model iterations and results.

| Analysis Method | Iteration Results | | |
|---|---|---|---|
| | | Fishnet2 | Fishnet3 |
| | *Distance* | *Yielded Results* | *Yielded Results* |
| **High/Low Clustering (Getis-Ord General G) HTML** | 164m | * | * |
| | 300m | * | * |
| | 500m | * | * |
| | 750m | * | * |
| | 1km | Yes | Yes |
| | 2km | Yes | Yes |
| | 5km | Yes | Yes |
| | 10km | Yes | Yes |
| **Multi-Distance Spatial Cluster Analysis (Ripley's K-Function)** | N/A | Yes | Yes |
| **Spatial Autocorrelation (Moran's I) HTML** | 164m | * | * |
| | 300m | * | * |
| | 500m | * | * |
| | 750m | * | * |
| | 1km | Yes | Yes |
| | 2km | Yes | Yes |
| | 5km | Yes | Yes |
| | 10km | Yes | Yes |
| **Cluster and Outlier Analysis (Anselin Local Moran's I)** | 164m | * | * |
| | 300m | Yes | Yes |
| | 500m | Yes | Yes |
| | 750m | Yes | Yes |
| | 1km | Yes | Yes |
| | 2km | Yes | Yes |
| | 5km | Yes | Yes |
| | 10km | Yes | Yes |
| **Hotspot Analysis (Getis-Ord Gi*)** | 164m | * | * |
| | 300m | Yes | Yes |
| | 500m | Yes | Yes |
| | 750m | Yes | Yes |
| | 1km | Yes | Yes |
| | 2km | Yes | Yes |
| | 5km | Yes | Yes |
| | 10km | Yes | Yes |

Table 3-1, continued:  Guide to model iterations and results.

There are a total of seven feature classes analyzed in this study (Table 3-1).  The

foreclosure points for Los Angeles County are utilized as the 'base layer' reference point

for aggregations and cluster analyses. The Assessor's Parcel Number, or APN, as well as spatial location for these parcels matched the provided foreclosure point dataset perfectly. No parcels were dropped from either dataset upon joining.

Table 3-1 describes the seven feature classes, the five cluster analysis methodologies, and the eight selected distance thresholds. Each possible permutation of these factors is listed in this table. Table 3-1 is intended to highlight which combinations of factors yielded results. Some of the results could not run due either to lack of sufficient computational power or a distance threshold which was too low (not accepted by the software program). Some of the results also were not applicable due to, for example, the foreclosure point dataset not being able to run the spatial autocorrelation tool against the foreclosure point dataset. The multi-distance spatial cluster analysis (Ripley's K-function) tool uses its own distance threshold variables and iteration counts when generating its functions (Cressie and Collins, 2001; Gatrell et al., 1996; Lentz, Blackburn and Curtis, 2011; Mitchell, 2009; Perry, Miller and Enright, 2006; Ripley, 1977). The results of Table 3-1 indicate only whether or not a given process ran successfully.

### 3.3 Data Preparation

To begin the dataset preparation process, parcels underwent the Dissolve process, where each polygon was 'blended' with each adjacent parcel which shared its boundaries. This new dataset is designated "Parcels2". Another aggregation process was performed

on Parcels2 to combine adjacent parcels within 5 meters from one another.  This

dissolution of boundaries created another new dataset referred to as "Parcels3".

A regular grid of polygons, 300 meters by 300 meters, was created using the

Create Fishnet tool in Esri's ArcGIS Toolbox.  This tool generates a regular grid of

rectangles across a given input, in this instance the extent of Los Angeles County.  This

new dataset is designated as "Fishnet1" and is comprised of 308,730 features.  Another

aggregation process was performed on this dataset to create "Fishnet2" with broader

dimensions, 1 kilometer by 1 kilometer, comprised of 27,798 features.  The final dataset

is referred to as "Fishnet3" and is identical in dimensions to Fishnet2, except that all

features are shifted by 500 meters along the horizontal axis to the East.  This shift is

intended to illustrate how precarious cluster analysis is on aggregated features when

boundaries are even slightly altered.  This shift is indicative of the zoning problem in the

MAUP.  Table 3-2 provides a summary of these datasets and how they are derived from

one another.

| Feature Dataset | Original Dataset | Derivation Procedure | Dimensions | No. of Features |
|---|---|---|---|---|
| FC_Points | Yes | Original Dataset | Within LA County | 55,631 |
| Parcels | Yes | Original Dataset | LA County | 2,367,742 |
| Parcel2 | No | Dissolved 'Parcels' | LA County | 84,907 |
| Parcels3 | No | Aggregated 'Parcels2' | LA County | 74,714 |
| Fishnet1 | Yes | Original Dataset | Grid covering LA County | 308,730 |
| Fishnet2 | No | Aggregated 'Fishnet1' | Grid covering LA County | 27,798 |
| Fishnet3_Shifted | No | Shifted features from 'Fishnet2' | Grid covering LA County | 27,798 |

Table 3-2:  Summary of the seven datasets utilized in this analysis.

**3.4 Cluster Analysis Methods**

Once all of the preparations for each dataset were completed, each dataset and its associated tool were placed into an ArcGIS model in Model Builder. Model Builder allows for multiple iterations of a single tool, or multiple tools, to run in sequence. A total of 137 model iterations were run which covered each dataset, distance threshold, and cluster analysis method, the results of which are shown in Table 3-1. Another example, Figure 3-1, highlights five model iterations within a single model in Model Builder which used the same input, but with different distance thresholds chosen. These models facilitated an efficient processing environment which also provided a convenient testing environment for the cluster tools and their outputs. Each feature class output was stored in a file geodatabase except those tools with tabular or graphic output. All tools used in this procedure are found in the Arc Toolbox from ESRI's ArcGIS software suite. Processing time varied across the different models and iterations. Z-scores were applied to the cluster analysis (hotspot and cluster and outlier) outputs which yielded datasets. The Z-scores for these datasets are used to generate maps of various areas of Los Angeles County to highlight differences between the outputs. Depictions of these results tables, graphs, datasets as well as the maps derived from them are included in Chapters Four and Five along with information on Z-scores and P-values.
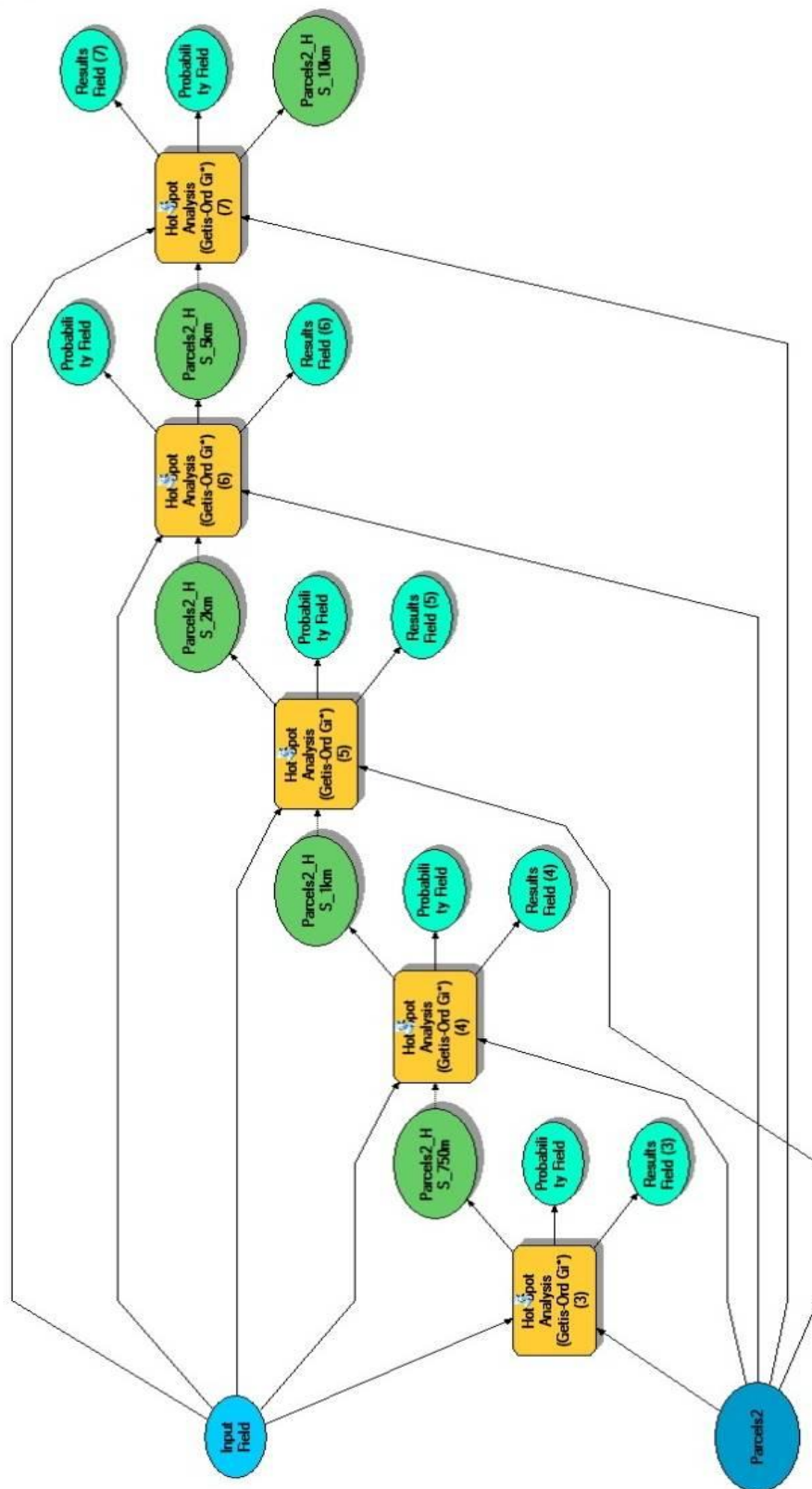
Figure 3-1: Example of a ModelBuilder model showing five Hotspot Analysis tools running in sequence. Yellow boxes represent the tools, dark blue ovals are dataset inputs, aqua colored ovals show the two result fields (Z-score and P-value), the aqua oval represents varying distance threshold inputs and the dark green ovals show the output datasets.

32

Many of the goals of this project involve illustrating areas where issues of either aggregation or zonal configuration can be clearly shown to switch between hot and cold spots. "Hot spots" delineate areas where the concentrations of a given variable (foreclosures) occur at higher concentrations. Conversely, 'cold spots' indicate areas where there is a marked absence of such occurrences. It must be noted that the red hot spots and blue cold spots of the Hot Spot Analysis tool differs in definition from the Cluster and Outlier Analysis tool. Figure 3-2 illustrates an example of the Hot Spot Analysis phenomenon wherein the left side of the figure shows that a series of proximal points is located within the same cell, thus that cell is highlighted in a red hue to indicate high clustering. The right area of the figure illustrates the same points, but dissected twice near the center of the cluster producing four equally sized cells with a similar number of points in each cell. The result of this shift is exemplified in the 500 meter shift between Fishnet2 and Fishnet3. This shift is meant to emphasize the zoning effect of the MAUP on the aggregated datasets. Since Fishnet2 is a regular grid that covers Los Angeles County indiscriminately of physical boundaries (such as lakes, freeways, and mountains) the shift of the units will highlight the zoning effect of the MAUP.
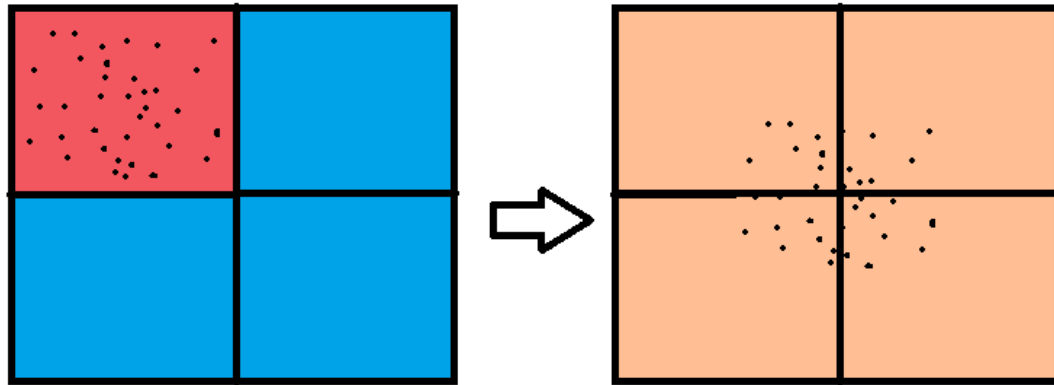
Figure 3-2: Identical sets of points that are either contained in a single cell (left) which generates a hotspot, or split into quadrants with different zonal boundaries (right) which shows even distribution.

## 3.5 Justification

The five different commonly used cluster analysis methods were chosen because such tools can visually illustrate how pervasive the MAUP is in various facets of spatial research in a statistically significant way. Each tool uses different functions to transform and analyze the spatial properties of different features within the same dataset. The wide variety of tools, aggregations, and distance thresholds lends credence to any possible conclusions regarding the MAUP impacts on data visualization in a spatial context. These permutations also exemplify how bias can play a role on how data is chosen, processed, and evaluated at a fundamental level. Each of the following chapters will detail specific attributes of the statistical tools used in this research endeavor.

Each of the five different methods discussed in this chapter were selected because of their inherent similarities and contrasting differences. Each method constitutes an aspect of analyzing clusters of spatial events within the same geometries. Examining the

MAUP requires looking at the same data for uniformity of experimentation while maintaining the integrity of a given tool. It might be possible to analyze the MAUP using only one cluster analysis method, though any results might be biased due to the nature of the data, the extent of the study area, any presumptions a researcher may or may not possess, and 'natural' clusters of spatially-oriented data. The elimination of these factors requires a comprehensive analysis method that draws from multiple cluster analysis tools, across multiple scales and multiple aggregations. A cross analysis encompassing multiple cluster tools is one of the best methods available for easy comparison, therefore it is implemented in this study.

**Chapter Four – Analysis Results**

Determining whether or not a feature's similarity to its neighbors is random

requires statistical tools for validation.  Each of the five tools utilized in this study

measure clustering in a different way.  It is important to emphasize that this thesis does

not reflect on all the details regarding how each tool analyzes a given dataset nor the

specific statistical equations used to garner the results.  Rather the goals of this research

emphasize the potential impact of scale and aggregation effects on clustering

methodologies and their results.

This chapter describes and illustrates the results of many of the model processes

for each tool's iteration.  To describe each result would be onerous; therefore this chapter

is broken down to convey the results of each tool individually, while Chapter Five

provides the interpretation of the results.  Each tool is summarized in an individual

subsection with supporting figures, tables, maps, and tool outputs pertinent to each

analysis method.  Hotspot analysis (Getis-Ord General G) and Cluster and Outlier

Analysis (Anselin Local Moran's I) have similar mapped outputs grouped into hot (red)

and cold (blue) zones to reflect clustering, discussed in sections 4.1 and 4.2, respectively.

These tools function as local indicators of spatial association, or LISA, statistics.  LISA

statistics give individual results for each unit in a dataset while global clustering tools

look at the datasets as a whole to determine whether or not clustering occurs.  The

following three tools act as indicators of global clustering.  Spatial autocorrelation

(Global Moran's I) and High/Low Clustering (Getis-Ord General G) have graphic outputs

that are summarized in tables in sections 4.3 and 4.4, respectively.  Finally, Multi-

Distance Spatial Cluster Analysis (Ripley's K-function) graphic results are displayed in section 4.5. Each subsection gives more detail on the functions and use of each tool as well as any relevant limitations of the results.

Before each section details the results, Figure 4-1 will be used to show all the instances of foreclosures in 2006-2008 in the Los Angeles County extent. Visual examination of this figure already illustrates specific foreclosure clustering in the denser urban areas. This figure can also be used as a guide to measure foreclosure instances compared to sections 4.1 and 4.2 which have similarly detailed maps.
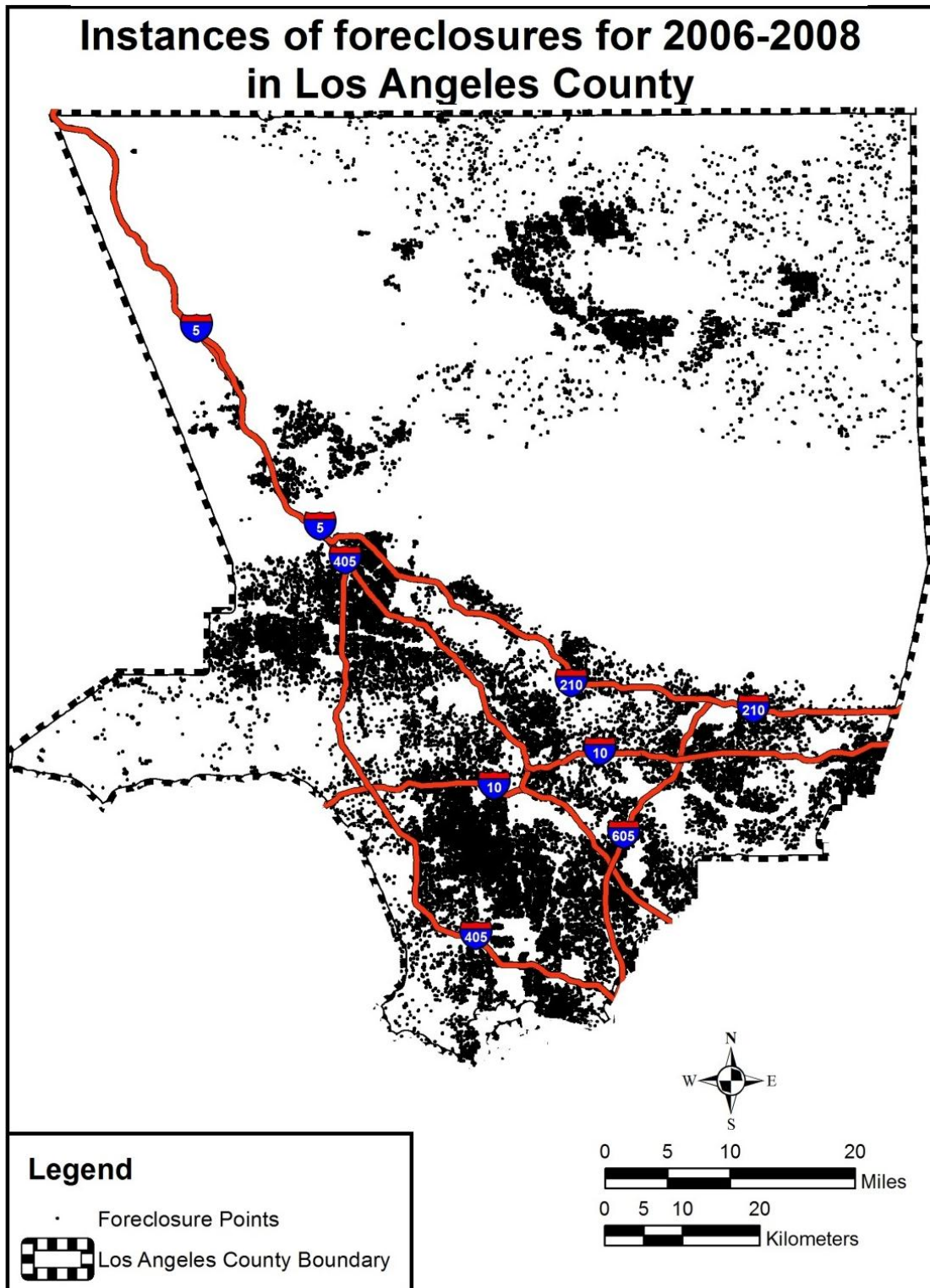
Figure 4-1: Overview map of the raw foreclosure instances in Los Angeles County from 2006 to 2008.

**4.1 Hotspot Analysis (Getis-Ord general G) Results**

Hotspot analysis is used to determine where spatial features cluster with respect to each other within a user-defined distance threshold (Ceccato, Haining, and Signoretta, 2002; Prendergast et al., 1993). This tool, as well as spatial autocorrelation and Cluster and Outlier analysis, are unable to process the foreclosure point dataset. Since hotspot analysis requires that each feature be weighted in some way, there was no method of running hotspot analysis on the foreclosure points without weighting them with some other unrelated variable. Such a weight would not allow the foreclosure points to be correlated with any of the other datasets whose weighted variable is a sum count of the total instances of foreclosure that take place within its boundaries. This exemplifies the aggregation effect of the MAUP. The foreclosure points cannot be compared to another variable without aggregating them to larger areal units. This justifies why the sum count of foreclosure instances within each polygon unit is used as the primary variable to be tested.

All of the other datasets were able to run at least five of the selected distance thresholds. The three parcels datasets ran successful iterations of hotspot analysis at a distance threshold of 750 meters and above. The three Fishnet datasets ran successfully at a distance threshold of 300 meters or more. Each iteration output resulted in a new dataset with new attribute fields pertaining to the p-value and the z-score of each feature's clustering properties. The null hypothesis of this analysis tool is that there is no inherent spatial association between the raw foreclosure points and any aggregations of them. P-values represent the probability that the null hypothesis was falsely rejected

(Mitchell, 2009).  Mitchell (2009) describes the Z-score as "a reference measure for the standard normal distribution (with mean of zero and standard deviation of 1)".  The Z-score is a useful way of determining whether or not to reject the null hypothesis.

Figure 4-2 provides an overview map of the cluster analysis for hotspots for the Los Angeles County extent.  Subsequent maps for this chapter are in a similar format. The dark red areas represent z-scores of 2 or higher, which indicate areas with high foreclosure instances, while the dark blue areas, with z-scores of -2 or less, showing areas contiguously void of foreclosures.  The light blue areas (z-score of 2.0 to 1.0) indicate areas that are moderately void of foreclosures while the light red or beige areas (z-score of -1.0 to -2.0) show areas with moderate foreclosure clustering.  The white areas represent areas of dispersion with neither high clustering nor low clustering.  In this analysis, 'dispersion' is defined as an area with neither a high concentration of foreclosures nor an absence of them.  Any green areas in the map are simply background imagery service showing forested areas via green polygons.  Figure 4-2 provides an overview map of the Parcels3 dataset, 2 kilometer distance threshold, with the contiguous Los Angeles County as the same scale.
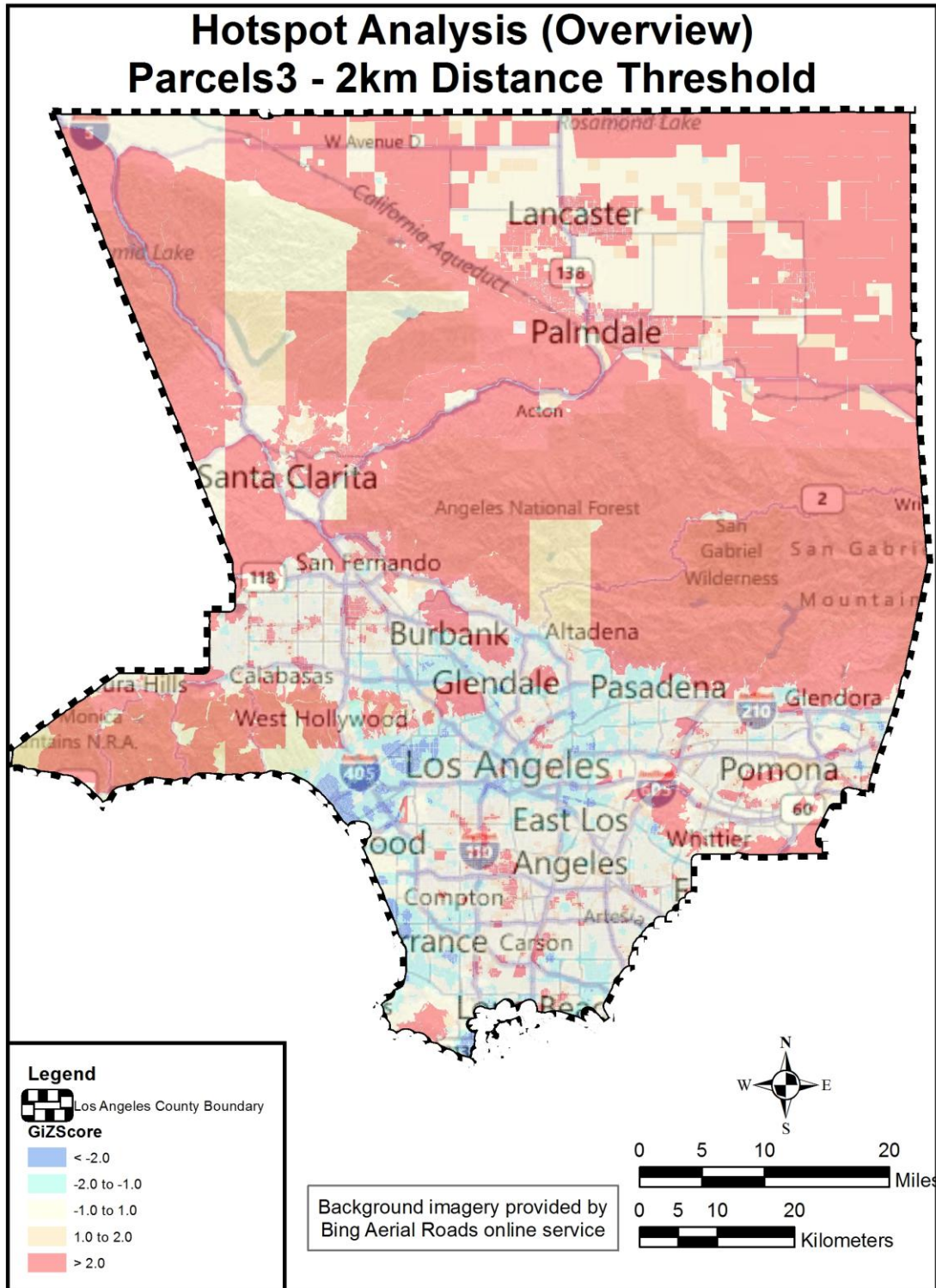
Figure 4-2: Overview hotspot analysis map of Parcels3 showing all Los Angeles County with a 2 kilometer distance threshold
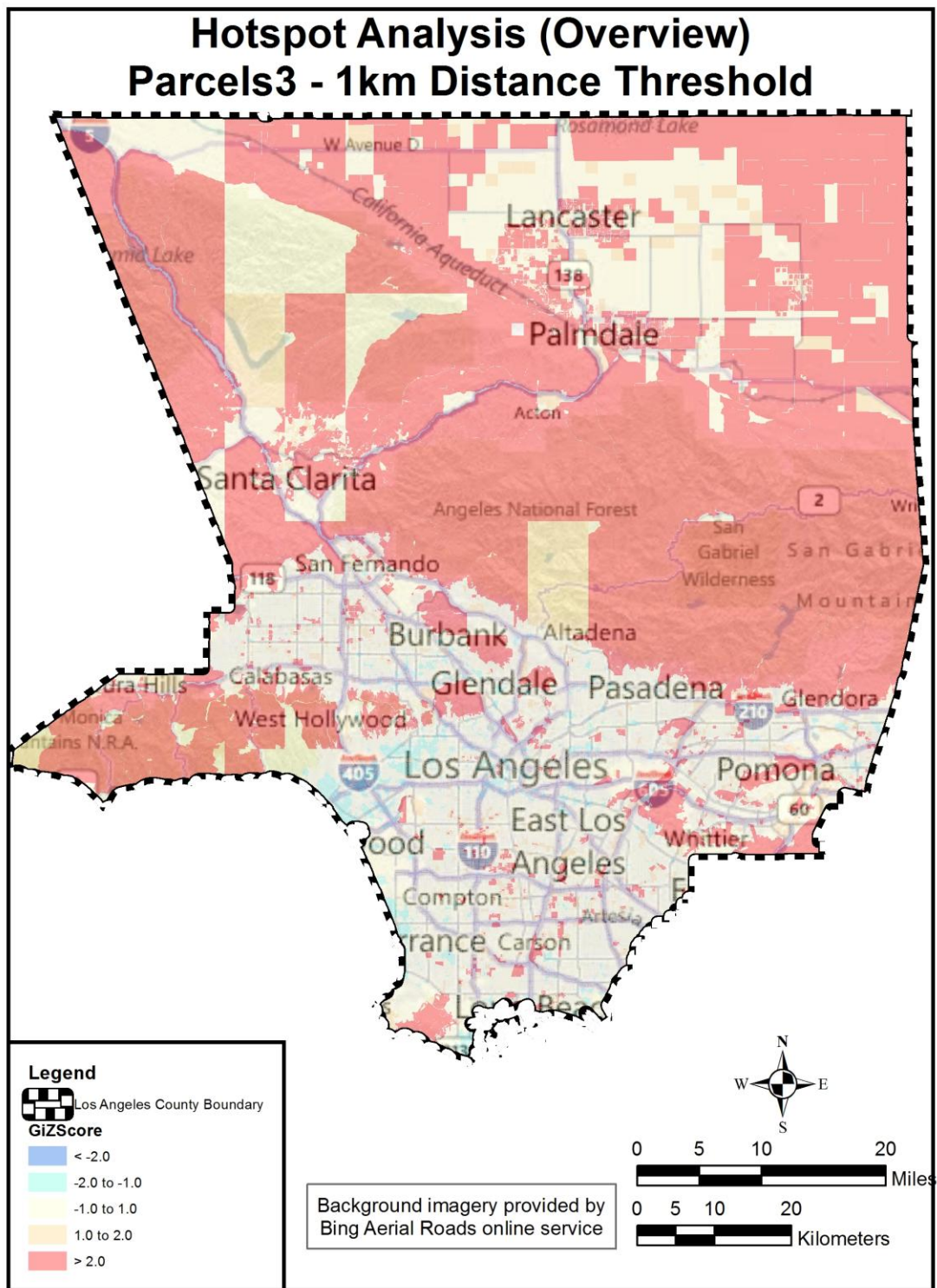
Figure 4-3: Overview hotspot analysis map of Parcels3 showing all Los Angeles County with a 1 kilometer distance threshold.

At smaller scales, some results/maps illustrate interesting differences when the distance threshold is shifted from 750 meters to 5 kilometers. Parcels2 results indicate severe shifts in hot/cold spot locations based on the distance threshold. Figure 4-4 and Figure 4-5 show the Parcels2 dataset zoomed into a smaller area to illustrate these shifts, how hotspot areas dramatically change when features are aggregated into larger units. This discussion and the following sections of this chapter are intended to highlight the final results of model iterations. Interpretation and further discussion of these results are provided in Chapter Five.
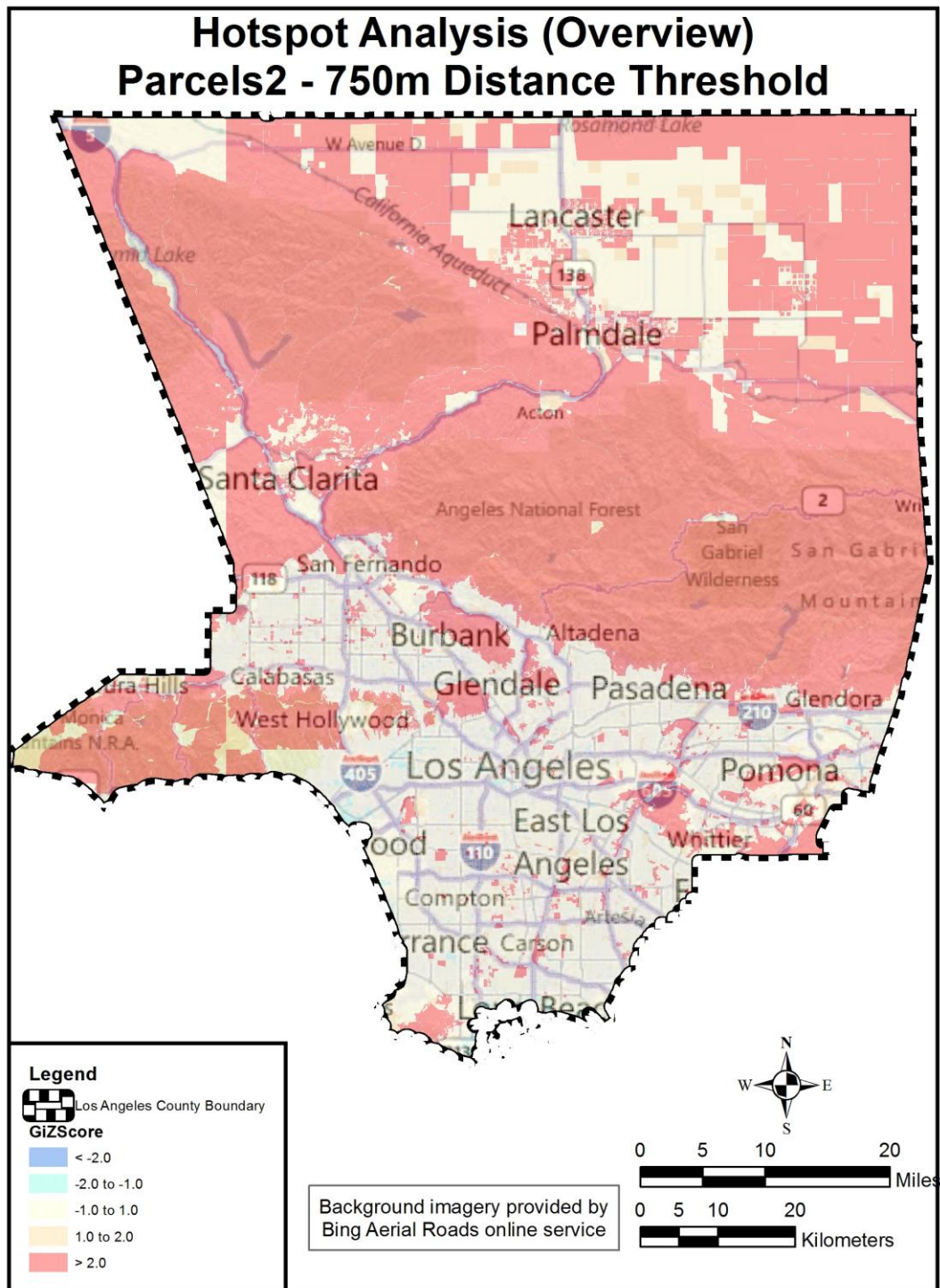
Figure 4-4:  Hotspot analysis results for the Parcels2 dataset at a distance threshold of 750 meters.

# Hotspot Analysis (Close View)
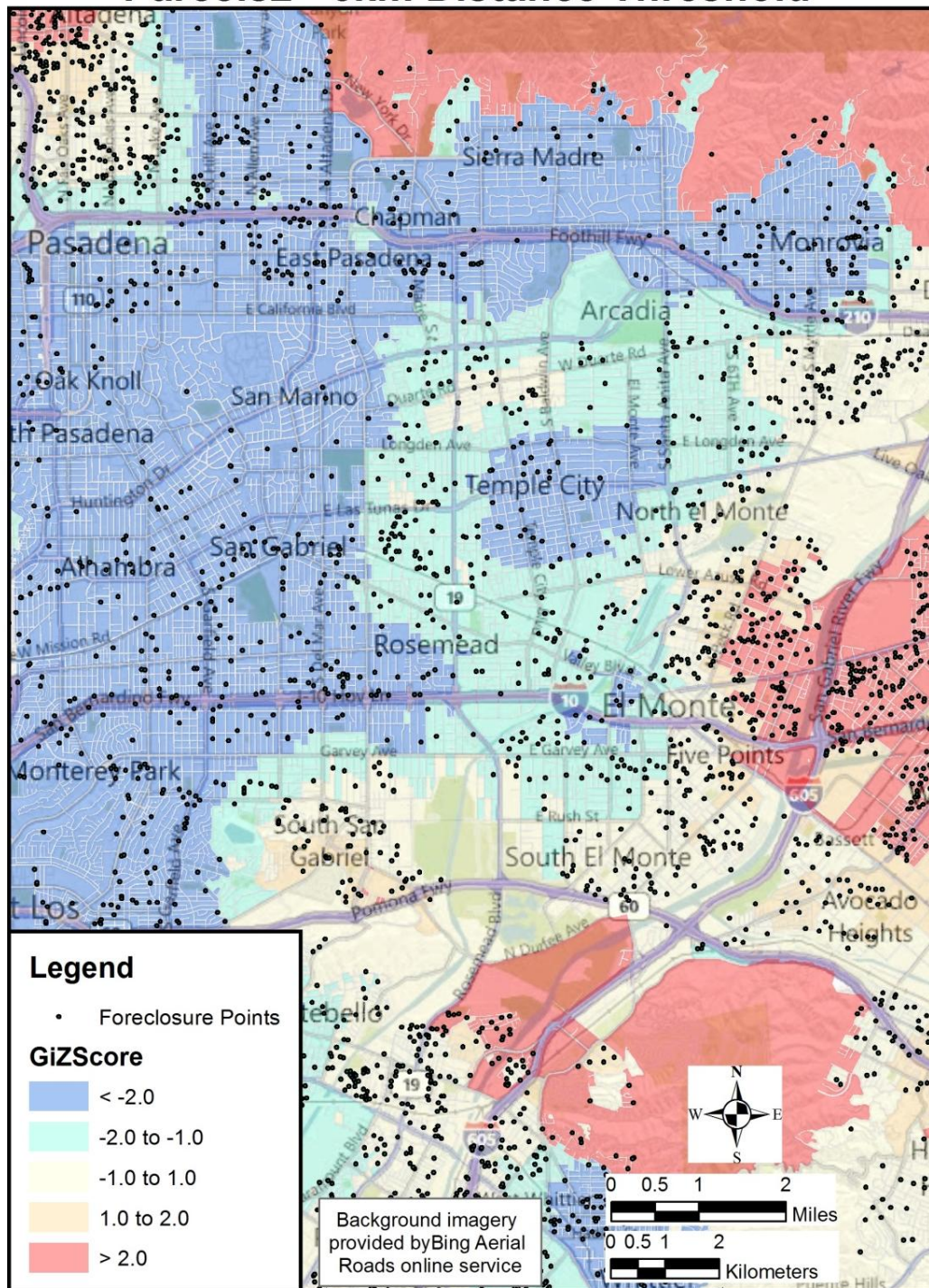## Parcels2 - 5km Distance Threshold



Figure 4-5:  Hotspot analysis results for the Parcels2 dataset at a distance threshold of 5 kilometers.

**4.2 Cluster and Outlier Analysis (Anselin Local Moran's I) Results**

      A cluster and outlier analysis is marked not necessarily by how individual features rate when clustered near one another, but rather by how their weights impact their standing in relation to all other features they are adjacent to (Mitchell, 2009). The results from an Anselin Local Moran's I test give a more detailed look at how features cluster to each other depending on the weighted attribute being tested against and the same weighted attributes from nearby features. An example of this basic principle can be found in Figure 4-6, where the Parcels2 dataset was processed with a 5 km distance threshold.

      It is important to reiterate that this analysis does not focus on explaining the differences between different permutations of model runs when changing certain aspects of how an analysis tool processes a given feature with its neighbors, but rather that almost all possible parameters remain the same with only the dataset and distance threshold changing. Altering these parameters in this tool or in any of the other related tools with similar parameters will result in more potential permutations than would be prudent to address, given the scope of this project.

      The outputs from the cluster and outlier analysis are indicative of spatial patterns that can be seen across multiple scales and aggregations. For example, Figure 4-7, which illustrates the results of the Fishnet1 dataset at a 5 kilometer distance threshold, shows clearly defined natural clusters at the county-wide scale. Conversely, when the dataset is shifted to Fishnet2, as in Figure 4-8, the same perspective or scale with the same distance

**Cluster & Outlier Analysis (Overview)**
**Parcels2 - 5km Distance Threshold**

Legend
Los Angeles County Boundary
**GiZScore**
- < -2.0
- -2.0 to -1.0
- -1.0 to 1.0
- 1.0 to 2.0
- > 2.0

Background imagery
provided byBing Aerial
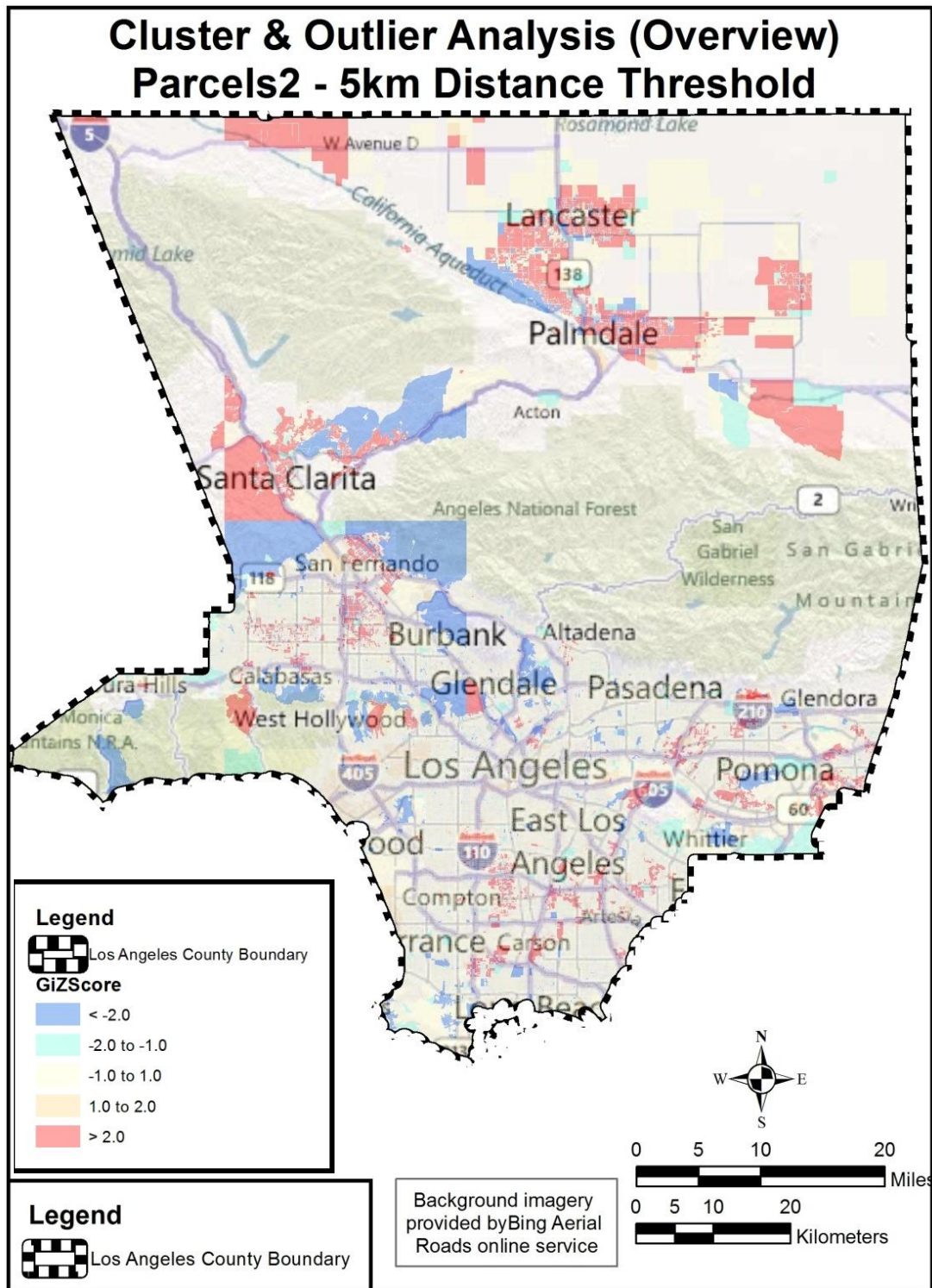Roads online service

Legend
Los Angeles County Boundary

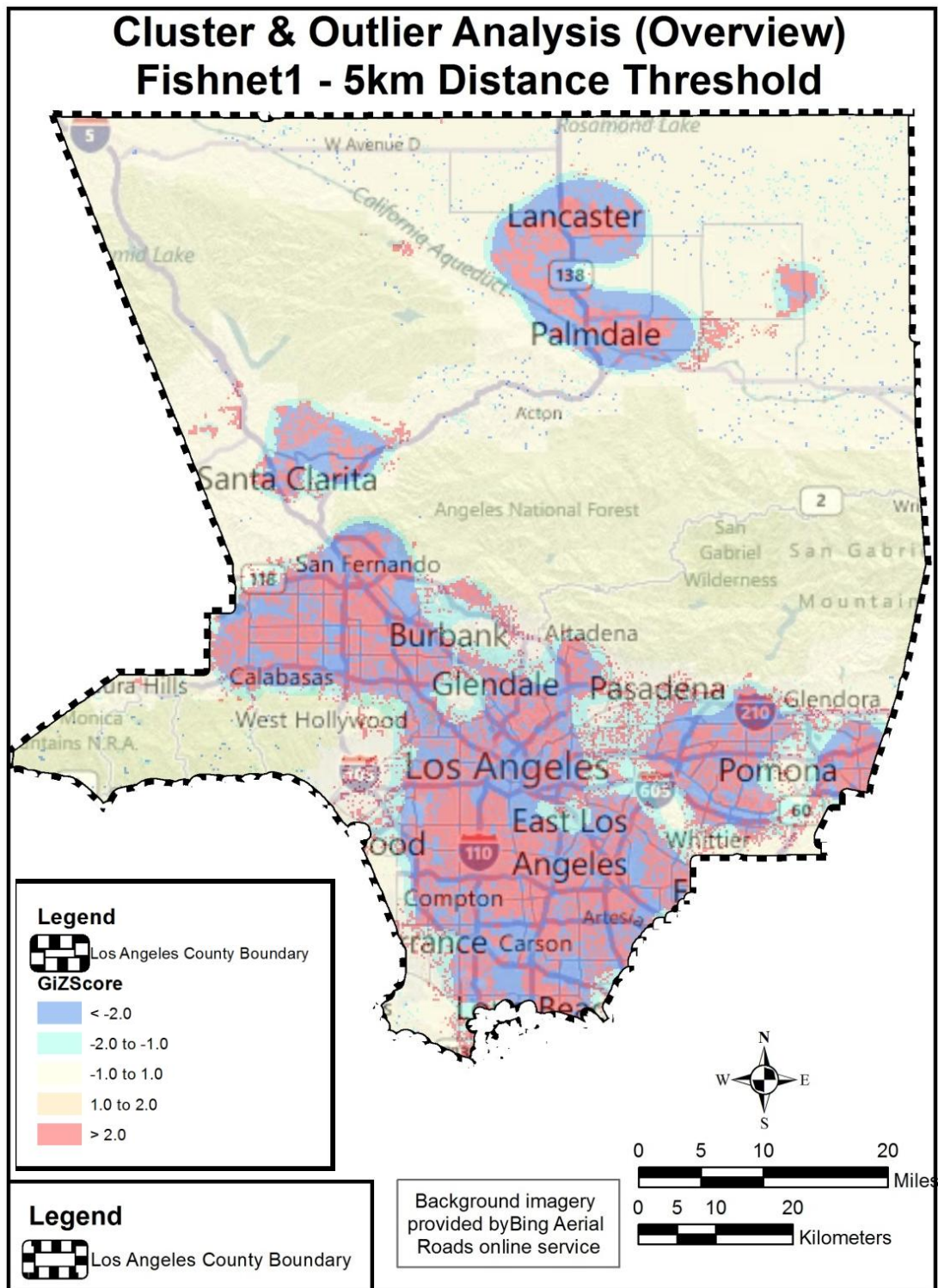Figure 4-6: Cluster and outlier analysis of the Parcels2 dataset with a 5km distance threshold.

Figure 4-7: Cluster and outlier analysis of the Fishnet1 dataset with a 5km distance threshold.
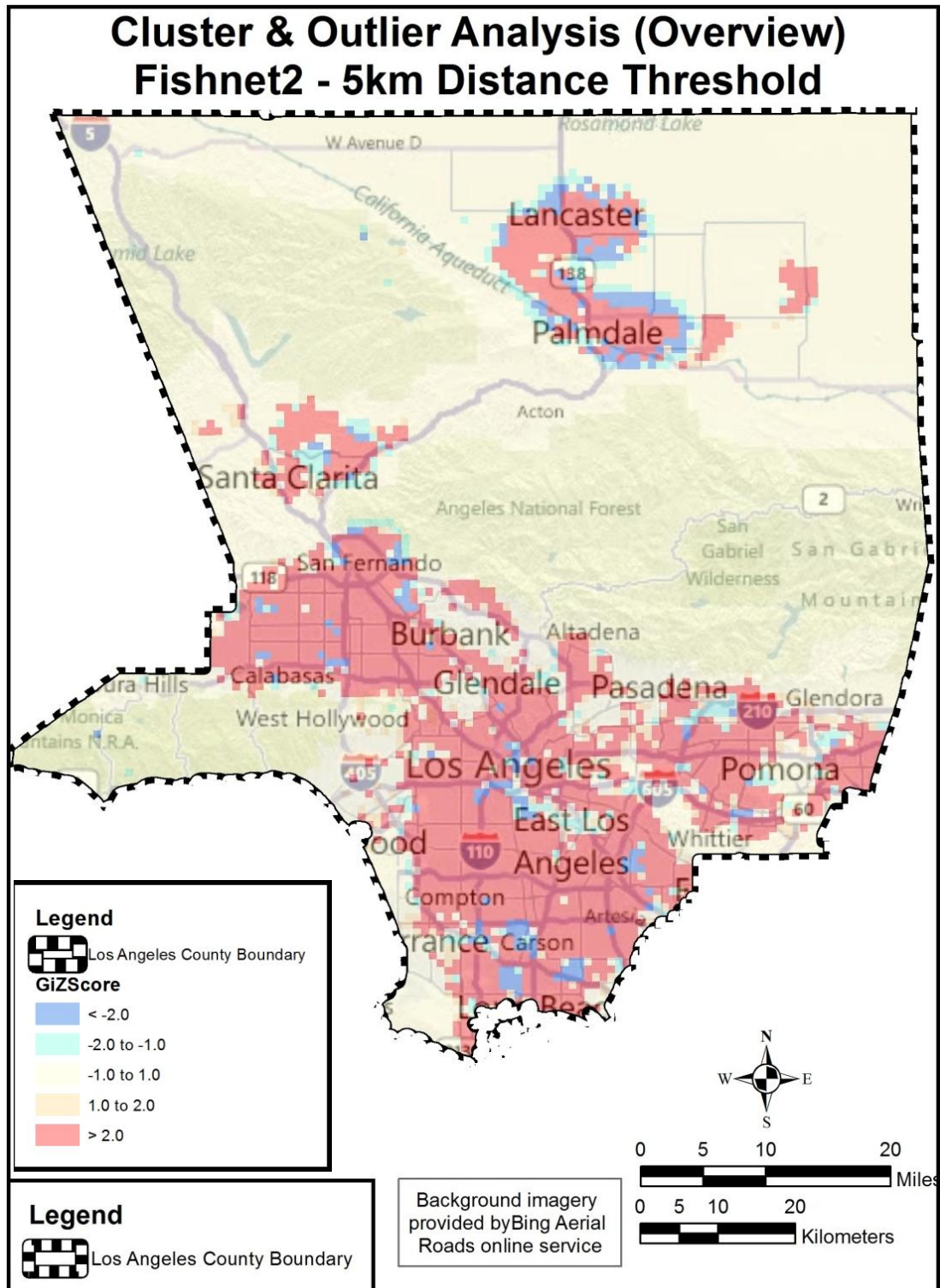
Figure 4-8: Cluster and outlier analysis of the Fishnet2 dataset with a 5km distance threshold.

threshold illustrates vastly different results.  There are now significantly higher numbers of contiguous hotspots, with a few cold spots that are difficult to discern.

One must remember that while the analysis tool utilizes similar outputs to the hotspot analysis tool (Getis-Ord General G), the interpretation of hot and cold spots is different.  This analysis looks at how similar a weighted feature is to every other feature within its immediate vicinity while taking into account all the other features within its distance threshold (Mitchell, 2009).  Areas that are 'cold', or dark blue, do not necessarily indicate areas where no foreclosures occur, but merely that the pattern of clustering varies greatly enough in that area to warrant a negative z-score and negative Moran's I value.  Alternately the hotspot analysis (Getis-Ord General G) determines whether or not a unit is 'hot' or 'cold' based on the total number and respective weights of all its neighbor features within the specified distance threshold.

## 4.3 Spatial Autocorrelation (Global Moran's I) Results

The Global Moran's I tool simultaneously measures a feature's weighted variable and its location to determine whether or not a pattern is dispersed, random, or clustered. It then assigns P-values and Z-scores to show the significance of the results (Dale et al., 2002; Mitchell, 2009).  All of the results for the Global Moran's I statistic indicated a positive likelihood that clusters were present in discernible patterns.  There was a less than 1% chance that patterns of clusters could have been the result of random chance. All Moran's Index scores are supposed to be between -1 and 1 (Mitchell, 2009). Negative Moran's I scores represent suspected dissimilarity, with an index closer to -1

representing the most dissimilarity. Conversely, positive Moran's I scores represent

suspected clustering, with values closer to 1 representing high levels of clustering with

similar values (Mitchell, 2009). The 1 kilometer distance threshold is the only model for

the Parcels dataset that ran successfully, though its results will not be considered due to

falling outside the valid range of results, see Table 4-1. The reasons for this aberration

are unknown. The remaining results from the Spatial Autocorrelation analysis are

provided in Table 4-1.

The Spatial Autocorrelation tool requires that each feature be weighted in some

way so as to indicate autocorrelation between and among various features and their

attributes. All of the aggregated datasets (except the foreclosure dataset) employed the

total count of all instances of foreclosure as the variable to be tested. Since the

foreclosure points dataset could not be tested against itself, it was exempted from this

particular analysis. This dataset could have been autocorrelated with another unrelated

variable, such as unemployment rates, though such a correlation would have had no

meaning compared to the other datasets analyzed or to their results.

All of the Moran's I scores and associated z-scores indicated strong clustering

patterns within and amongst the various datasets. Each of these results corresponds with

the results from the Ripley's K-function results as shown in section 4.5. This

corroboration indicates not only that the clustering of foreclosures is indicative of

inherent assemblages of foreclosure instances but that their groupings tend toward

recognizable spatial patte

| Spatial Autocorrelation Results | | | |
|---|---|---|---|
| Dataset | Distance | Z-Score | Moran's Index |
| **Parcels** | 164m* | * | * |
| | 300m* | * | * |
| | 500m* | * | * |
| | 750m* | * | * |
| | 1km** | 105500.436** | 31.702** |
| | 2km | Incomplete | Incomplete |
| | 5km | Incomplete | Incomplete |
| | 10km | Incomplete | Incomplete |
| **Parcels2** | 164m* | * | * |
| | 300m* | * | * |
| | 500m* | * | * |
| | 750m | 32.534 | 0.062 |
| | 1km | 41.102 | 0.061 |
| | 2km | 71.654 | 0.06 |
| | 5km | 130.716 | 0.054 |
| | 10km | 193.249 | 0.048 |
| **Parcels3** | 164m* | * | * |
| | 300m* | * | * |
| | 500m* | * | * |
| | 750m | 33.617 | 0.077 |
| | 1km | 42.449 | 0.075 |
| | 2km | 74.593 | 0.074 |
| | 5km | 136.533 | 0.066 |
| | 10km | 202.423 | 0.058 |

Table 4-1:  Summary of  the Spatial Autocorrelation results.
*This particular dataset could not be run due to an insufficient distance threshold within the extent of LA County for the tool to run properly.
**Results do not fall within acceptable parameters.

| Spatial Autocorrelation Results | | | |
|---|---|---|---|
| **Fishnet1** | 164m* | * | * |
| | 300m* | * | * |
| | 500m | 562.288 | 0.514 |
| | 750m | 789.443 | 0.476 |
| | 1km | 972.212 | 0.451 |
| | 2km | 1516.324 | 0.392 |
| | 5km | 2587.939 | 0.297 |
| | 10km | 3564.576 | 0.225 |
| **Fishnet2** | 164m* | * | * |
| | 300m* | * | * |
| | 500m* | * | * |
| | 750m* | * | * |
| | 1km | 176.278 | 0.749 |
| | 2km | 270.711 | 0.691 |
| | 5km | 491.362 | 0.539 |
| | 10km | 682.698 | 0.414 |
| **Fishnet3** | 164m* | * | * |
| | 300m* | * | * |
| | 500m* | * | * |
| | 750m* | * | * |
| | 1km | 174 | 0.739 |
| | 2km | 267.137 | 0.681 |
| | 5km | 485.489 | 0.532 |
| | 10km | 674.508 | 0.409 |

Table 4-1, continued:  Summary of Spatial Autocorrelation results.
*This particular dataset could not be run due to an insufficient distance threshold within the extent of LA County for the tool to run properly.

## 4.4 High/Low Clustering Results (Getis-Ord General G) Results

The High/Low Clustering tool (Getis-Ord General G) calculates a single statistic

for a given dataset at different distance thresholds, one statistic per iteration, rather than

statistical results for each feature.  This tool provides the benefit of describing whether or

not a dataset indicates if there are concentrations of high or low values and whether or not these values are clustered. A tool such as this works best in describing the overall trend of a given study area, depending on the scope of one's research and goals (Mitchell, 2009). The primary results of this analysis indicate that all successful iterations resulted in a statistically significant indication of clustering showing a p-value of less than 0.01, with the sole exception of Parcels2 at the 1 kilometer distance threshold, which had a p-value of 0.06.

The results of the high/low clustering analysis match well with the results of the spatial autocorrelation analysis. Global Moran's I and Getis-Ord general G both measure continuous feature values in order to ascertain an overall pattern of the dataset. These results indicate that the local level of statistics in the results of other tools' processes should provide a clearer idea as to the spatial patterns involved in the instances of foreclosures across various scales, aggregations, and datasets.

| Summary of High/Low Clustering Results | | | | |
|---|---|---|---|---|
| *Dataset* | *Distance* | *Z-Score* | *Observed General G* | *p-value* |
| **Parcels** | 164m | * | - | - |
| | 300m | * | - | - |
| | 500m | * | - | - |
| | 750m | 36971.598 | 0.005192 | 0 |
| | 1km | 36964.607 | 0.005192 | 0 |
| | 2km | ** | - | - |
| | 5km | ** | - | - |
| | 10km | ** | - | - |

Table 4-2: Summary of High/Low Clustering results.
*This particular dataset could not be run due to an insufficient distance threshold for the tool to run properly.
**This dataset could not run due to insufficient processing power of available hardware.

| Summary of High/Low Clustering Results | | | | |
|---|---|---|---|---|
| | 164m | * | - | - |
| | 300m | * | - | - |
| | 500m | * | - | - |
| | 750m | 2.876 | 0 | 0.004029 |
| Parcels2 | 1km | 1.879 | 0 | 0.060241 |
| | 2km | 6.808 | 0.000001 | 0 |
| | 5km | 8.456 | 0.000002 | 0 |
| | 10km | ** | | |
| | 164m | * | | |
| | 300m | * | | |
| | 500m | * | | |
| | 750m | 9.233 | 0 | 0 |
| Parcels3 | 1km | 8.287 | 0 | 0 |
| | 2km | 15.717 | 0.000001 | 0 |
| | 5km | 19.691 | 0.000002 | 0 |
| | 10km | 17.507 | 0.000003 | 0 |
| | 164m | * | | |
| | 300m | 425.936 | 0.000001 | 0 |
| | 500m | 562.338 | 0.000001 | 0 |
| | 750m | 789.321 | 0.000002 | 0 |
| Fishnet1 | 1km | 971.59 | 0.000003 | 0 |
| | 2km | 1507.599 | 0.000006 | 0 |
| | 5km | 2407.209 | 0.000012 | 0 |
| | 10km | 2522.363 | 0.000019 | 0 |
| | 164m | * | | |
| | 300m | * | | |
| | 500m | * | | |
| | 750m | * | | |
| Fishnet2 | 1km | 176.384 | 0.000002 | 0 |
| | 2km | 270.669 | 0.000004 | 0 |
| | 5km | 482.874 | 0.00001 | 0 |
| | 10km | 613.982 | 0.000017 | 0 |

Table 4-2, continued:  Summary of High/Low Clustering results.
*This particular dataset could not be run due to an insufficient distance threshold for the tool to run properly.
**This dataset could not run due to insufficient processing power of available hardware.

| Summary of High/Low Clustering Results | | | | |
|---|---|---|---|---|
| | 164m | * | - | - |
| | 300m | * | - | - |
| | 500m | * | - | - |
| Fishnet3_Shifted | 750m | * | - | - |
| | 1km | 174.107 | 0.000002 | 0 |
| | 2km | 267.102 | 0.000004 | 0 |
| | 5km | 477.224 | 0.00001 | 0 |
| | 10km | 607.364 | 0.000017 | 0 |

Table 4-2, continued:  Summary of High/Low Clustering results.
*This particular dataset could not be run due to an insufficient distance threshold for the tool to run properly.
**This dataset could not run due to insufficient processing power of available hardware.

## 4.5 Multi-Distance Spatial Cluster Analysis (Ripley's K-function) Results

Haase (1995) and Diggle et al. (2003) both note how the Ripley's K statistic tests against a null hypothesis which supposes that a distribution of points is completely random.  When the Ripley's K statistic is run on a dataset of spatial points it can be determined whether or not there is any inherent clustering of said points.  A procedure involving the production of random points within the same extent of the dataset and then running simulations to determine a 'confidence envelope' is described by Haase (1995). This confidence envelope is used to illustrate a 95%, 99%, or 99.9% confidence level of statistical accuracy based on the number of simulations run per dataset.  This project used a 99% confidence envelope for the five datasets which could be processed.

Unlike the other tools utilized in this project, the Multi-Distance Spatial Cluster Analysis tool has multiple distance thresholds built into the overall structure of the tool processes (Dale et al., 2002; Gatrell et al, 1996; Lentz, Blackburn and Curtis; 2011).  The

foreclosure points dataset has a much higher Observed K than the Expected K, which

indicates a strong clustering pattern, see Figure 4-9.  There is no statistically significant

difference between the expected and observed K-values for either of the two Fishnet

datasets capable of being calculated, see Figure 4-10 and Figure 4-10.  Fishnet2 and

Fishnet3 had similar results due to their almost identical features and extents (in terms of

total area and dimensions).  This is to be expected given the nature of the Fishnet

composition.  Without variations in the perimeter or area of each polygon, it will be

difficult for clustering to occur outside of an expected range.  Unfortunately, the Parcels

and Fishnet1 datasets were unable to be computed given the nature of the dataset and the

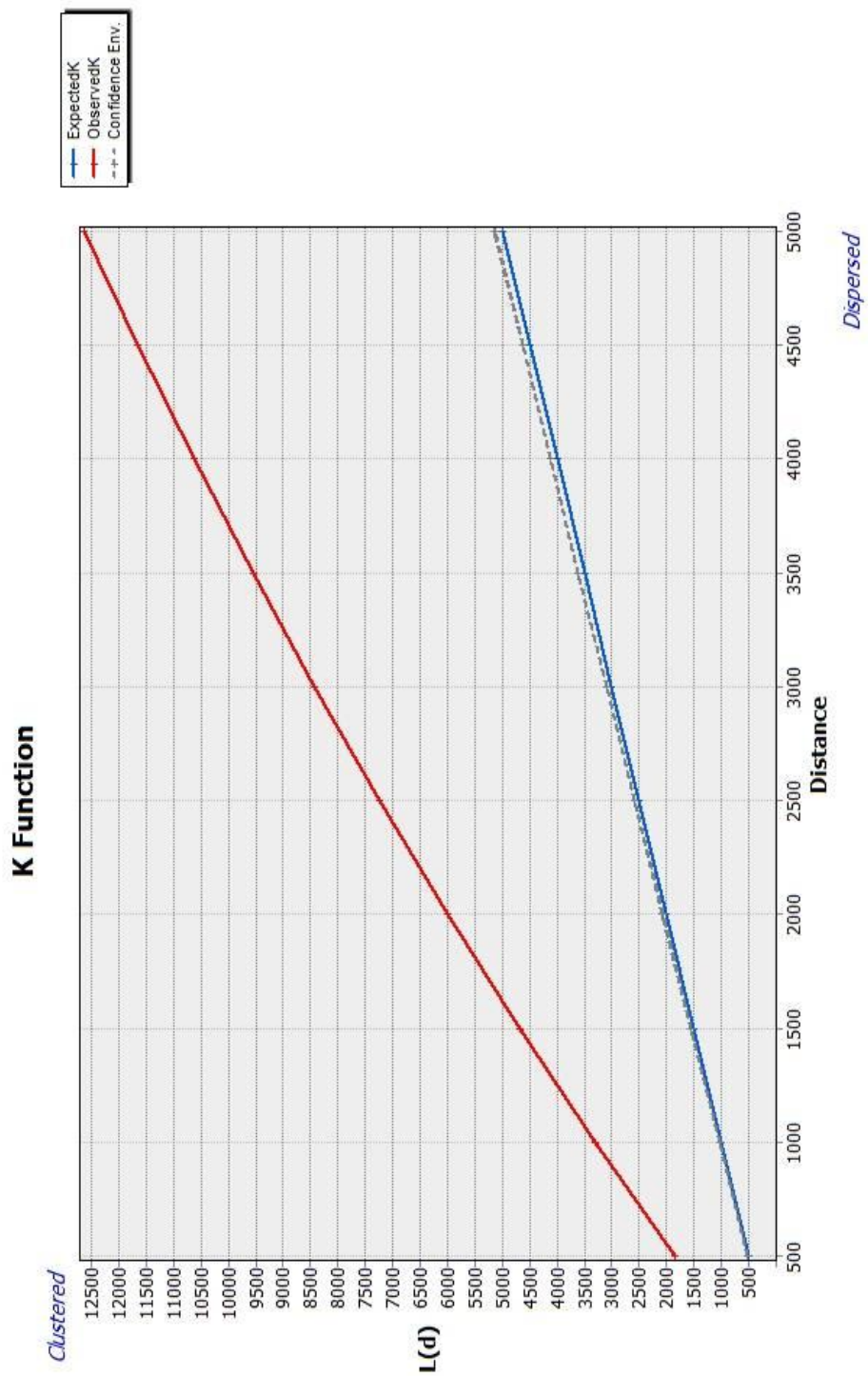processes involved in using this tool.

Figure 4-9: Ripley's K graph showing the Expected K, Observed K, and Confidence Interval for the Foreclosure Points dataset.
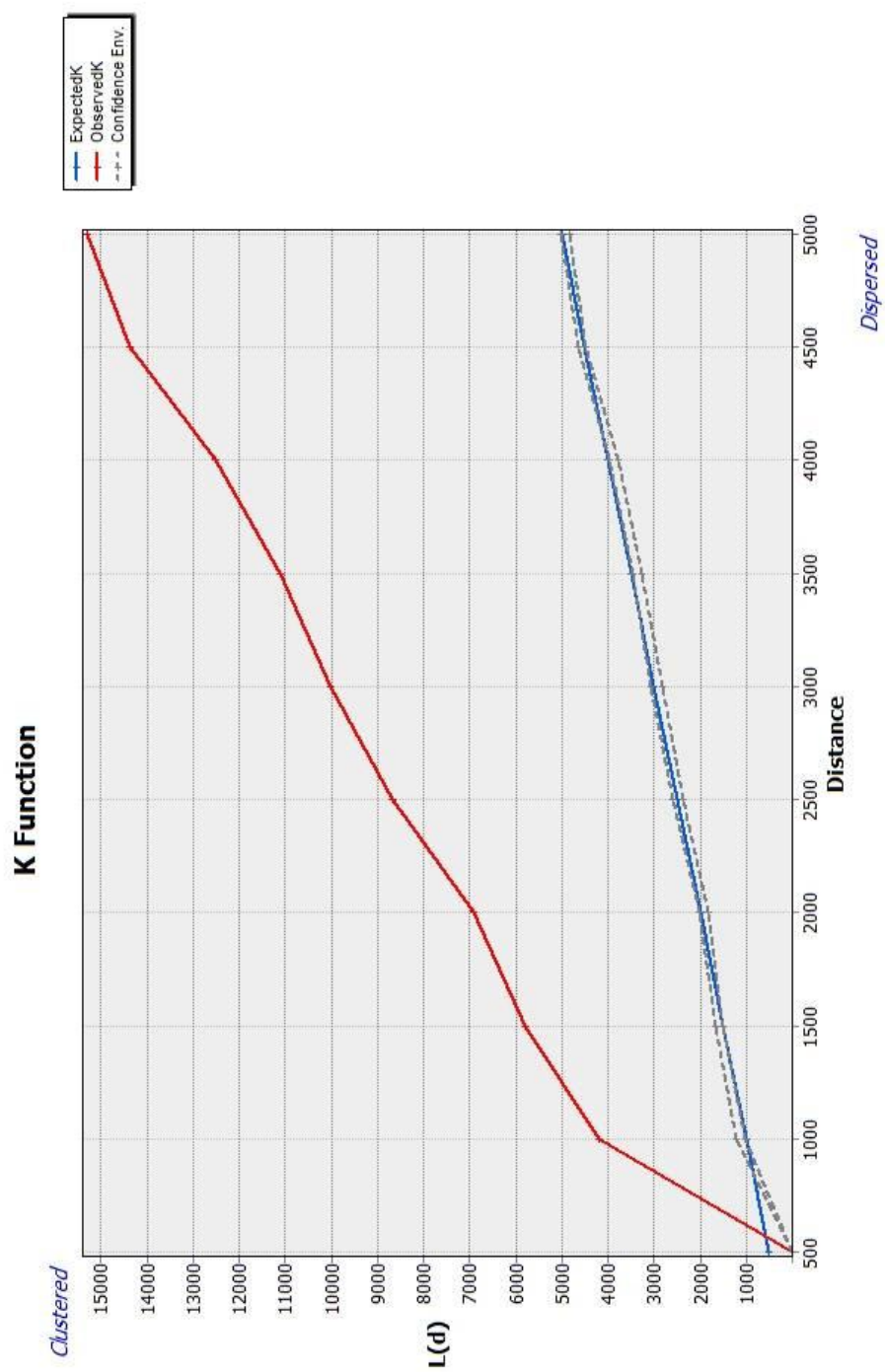
Figure 4-10: Ripley's K graph showing the Expected K, Observed K, and Confidence Interval for the Fishnet2 dataset.
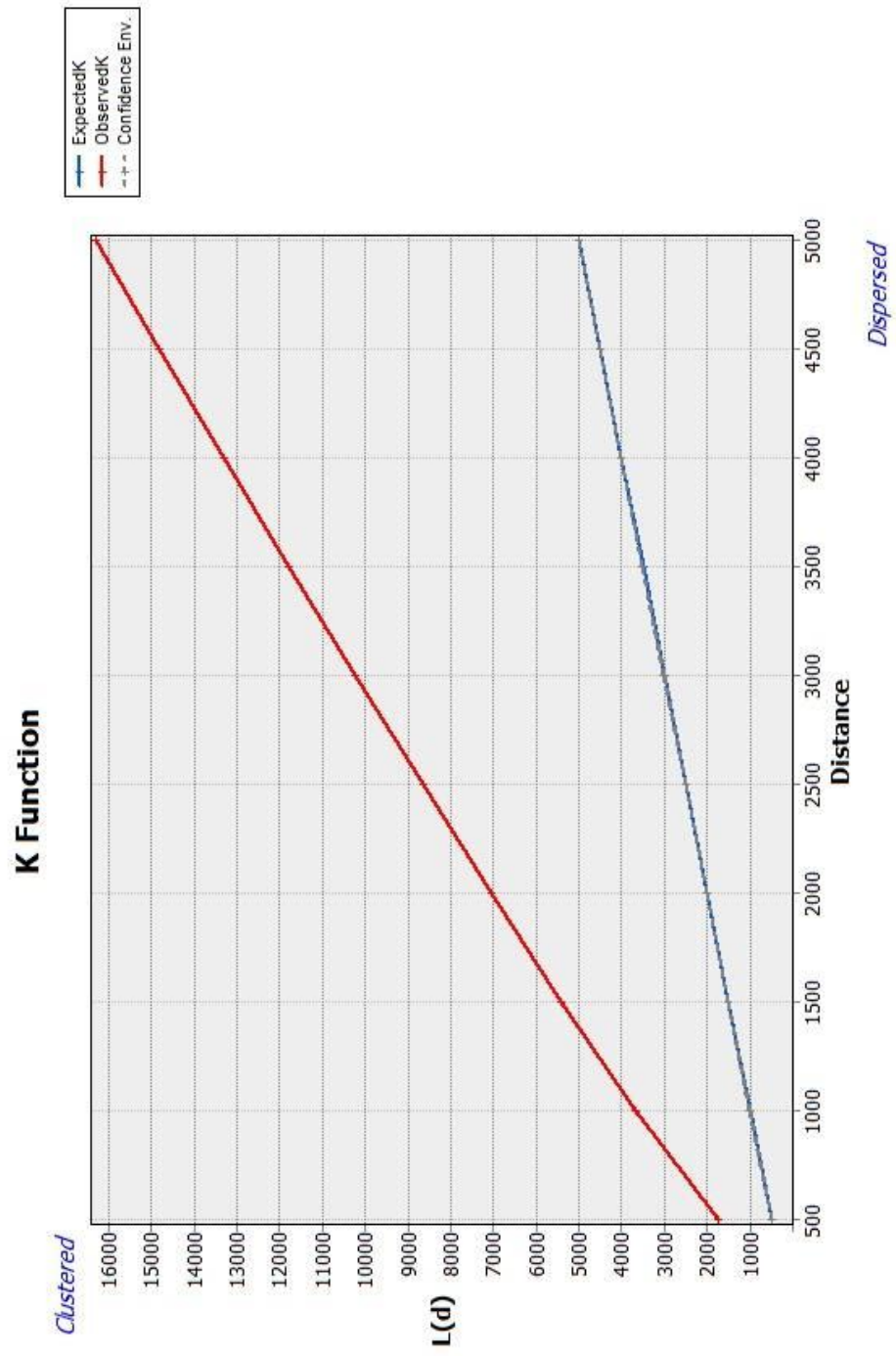
Figure 4-11: Ripley's K graph showing the Expected K, Observed K, and Confidence Interval for the Parcels2 dataset.
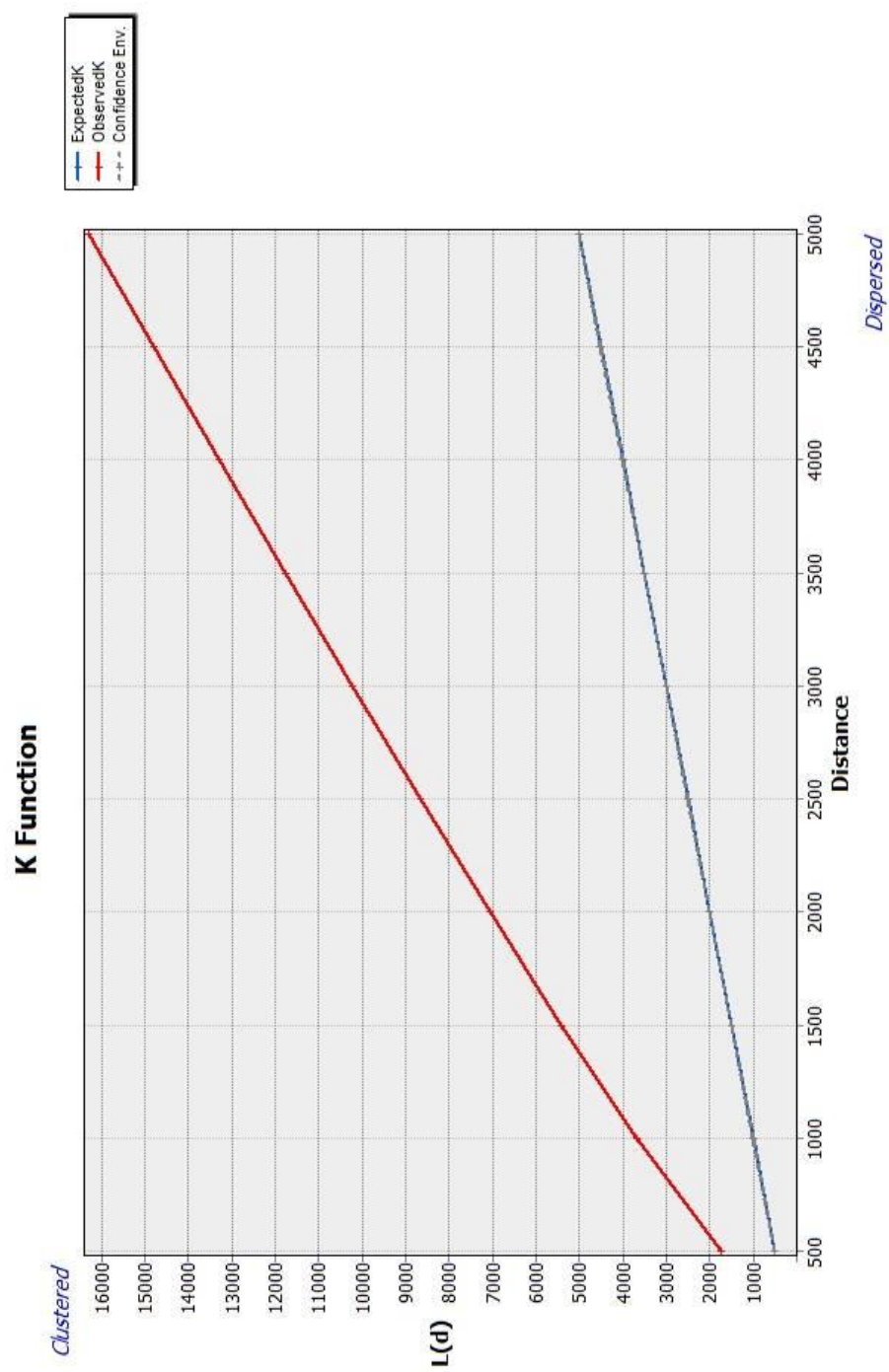
Figure 4-12: Ripley's K graph showing the Expected K, Observed K, and Confidence Interval for the Parcels3 dataset.

The Parcels2, Parcels3, Fishnet2, and Fishnet3 K-functions show a clear indication of strong clustering within the instances of foreclosure in Los Angeles County. These results prove that the null hypothesis is false which lends credence and validity to the other cluster analysis results. The Multi-Distance Spatial Cluster Analysis tool utilizes a beginning distance threshold (in this project, 500 meters is utilized) with a total of ten iterations per dataset, each with an additional 500 meter distance threshold. The latter is shown in Figure 4-9 to 4-13 as the independent variable on the X-axis. The perpetually increasing distance thresholds in this analysis, as well as the eight distance thresholds selected for the other tools, act as a gauge for how scale effects impact the way in which the results can be perceived as being clustered or dispersed, depending on the dataset and distance threshold chosen. Figures 4-10 and 4-11 illustrate this principle when the distance threshold is between 500 meters and 1 kilometer, which indicates that the Fishnet2 and Fishnet3 datasets have a high level of dispersion. All the distance thresholds at or above 1 kilometer show patterns of significant clustering. Future spatial research along similar veins must remain cognizant of how much of an impact scale has on the final results of any cluster analysis. Overall, the Multi-Distance Spatial Cluster Analysis illustrated the inherent clustering of the foreclosure point aggregations as well as other clustering, based on other datasets of varying aggregations, scales, and distance thresholds.
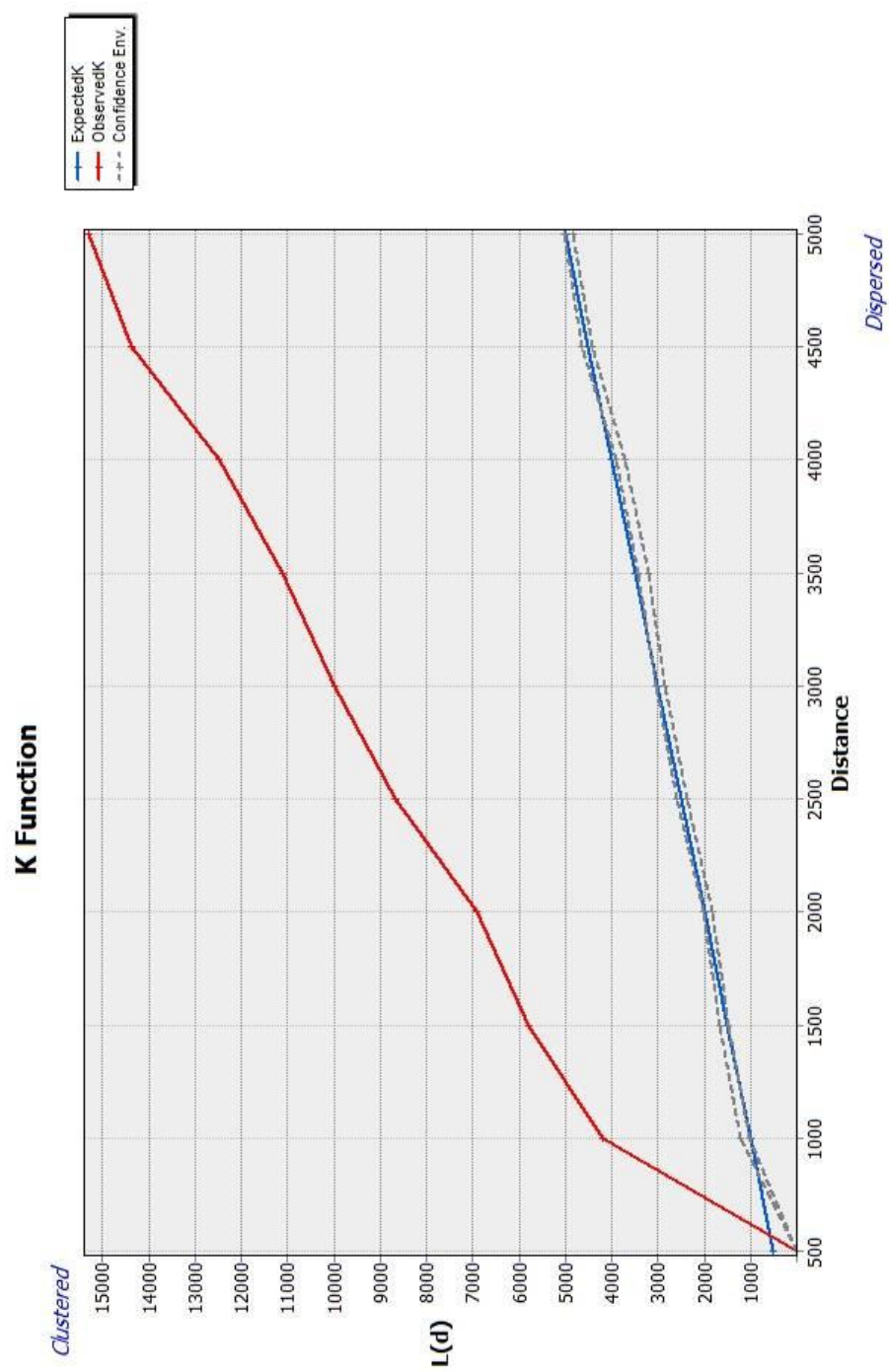
Fig 4-13: K-function graph showing the Expected K, Observed K, and Confidence Interval for the Fishnet3 dataset.

**Chapter Five – Interpretation**

**5.1 Interpretation of Hotspot Results (Getis-Ord Gi*)**

Much of Chapter Four dealt with assessing the individual results of each model iteration with summaries or examples. These results varied, as indicated in this chapter with tables, graphs, and maps. The emphasis of Chapter Five is to compare the results of the two LISA tools (Hotspot and Cluster and Outlier Analyses). These two sets of outputs are more easily compared to one another since they analyze features individually rather than entire datasets. The results are shown to exemplify the effects of the MAUP across multiple distance thresholds and areas. Datasets that share similar scales and aggregation methodologies provide a greater understanding of how the MAUP interacts with spatial data, analysis results, and interpretation(s) of those results.

Hotspot maps detailing a small area within Los Angeles County help provide a means for understanding just how prevalent the MAUP is in spatial research. Figures 5-1 and 5-2 show an area near Rosemead, California that highlights some important differences in hotspot clusters foreclosures. Figure 5-1 shows the hotspot Gi Z-score values on a map for the Parcels3 dataset at a 750 meter distance threshold. Figure 5-2 shows the same area with the same dataset but with a 10 kilometer distance threshold. Many of the areas in Figure 5-1 that were considered to be 'cold' spots are now shown to be dispersed in Figure 5-2.

# Hotspot Analysis (Close View)
## Parcels3 - 750m Distance Threshold



**Legend**

- Foreclosure Points

**GiZScore**

- < -2.0
- -2.0 to -1.0
- -1.0 to 1.0
- 1.0 to 2.0
- > 2.0

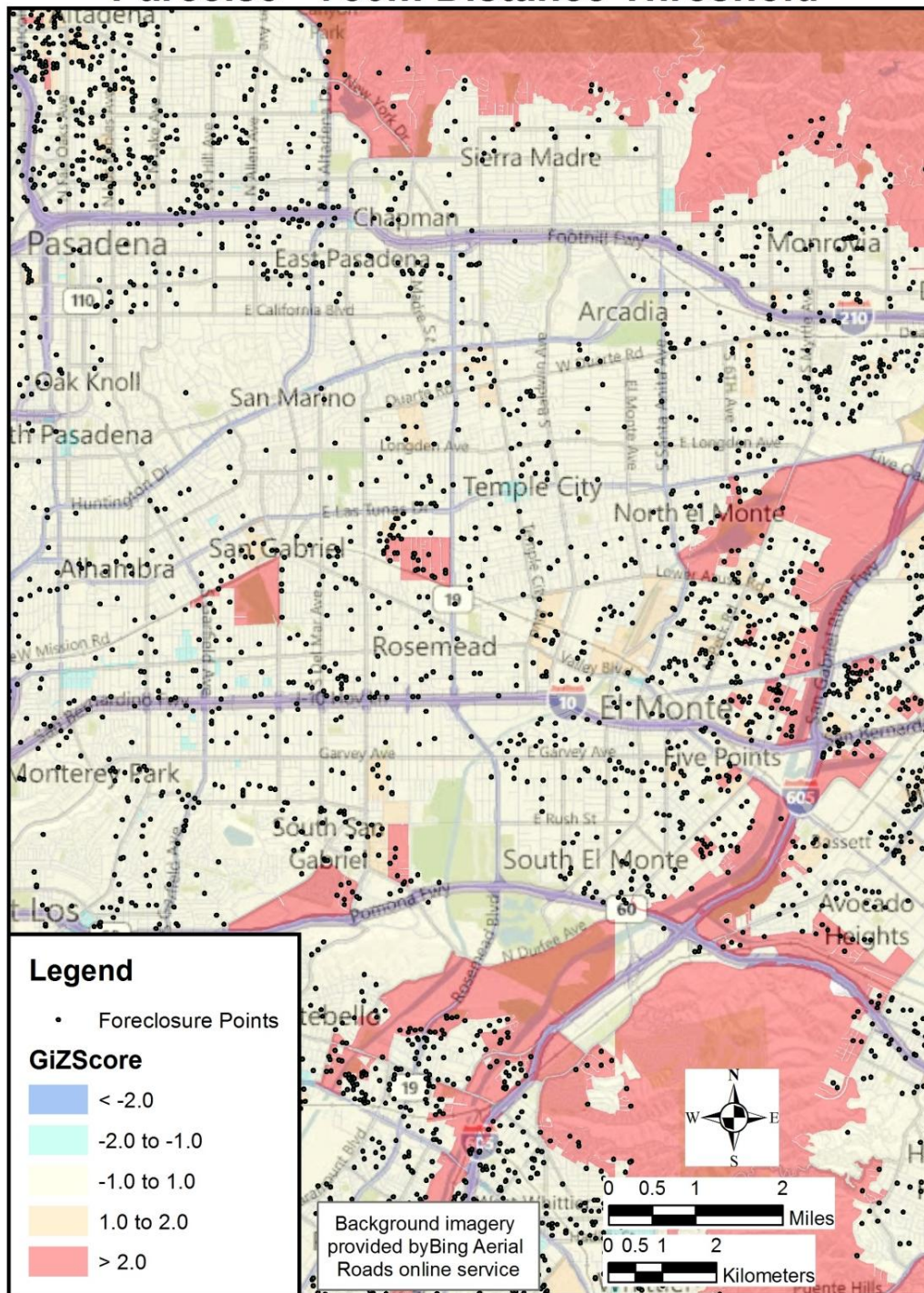Background imagery provided byBing Aerial Roads online service

Figure 5-1: Hot spot analysis map of Parcels3 dataset at a distance threshold of 750 meters.

# Hotspot Analysis (Close View)
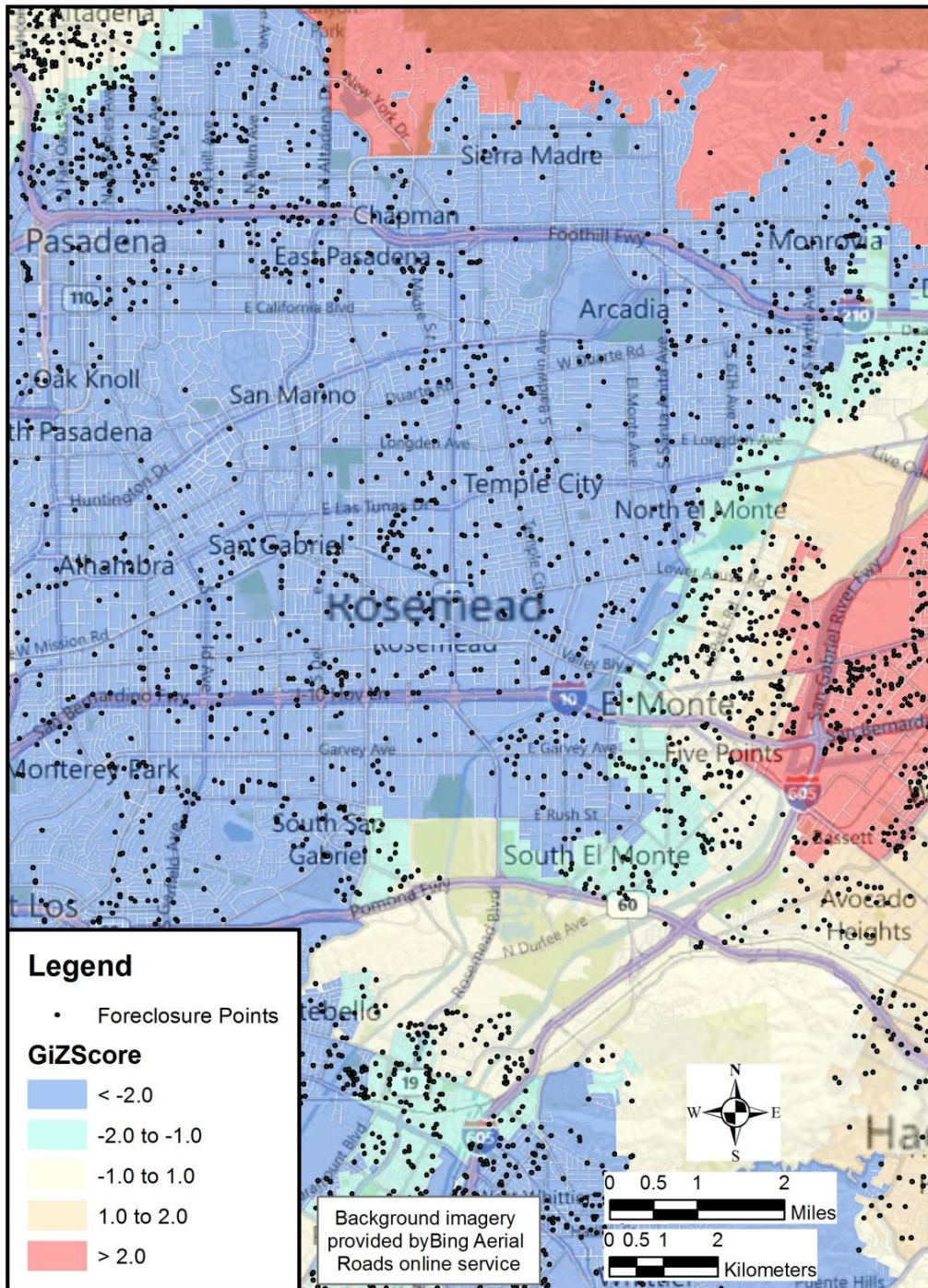# Parcels3 - 10km Distance Threshold



Figure 5-2: Hot spot analysis map of Parcels3 dataset at a distance threshold of 10 kilometers.

Particular interest should be paid to the red hotspot near the center of Figure 5-1. This area is shown with greater detail in Figures 5-3 and 5-4. Figure 5-3 shows the same dataset, Parcels3, with a distance threshold of 750 meters while Figure 5-4 shows the area with a 10 kilometer distance threshold. It is noteworthy how Figure 5-3 shows an area with a few hotspots, while the rest of the surroundings appear to have a dispersed number of foreclosures. This is indicated by the surrounding area having a GiZScore of -1.0 to 1.0 (beige coloring on the map). Figure 5-4, with its differing distance threshold, shows a far greater shift from the previous figure, Figure 5-3. Figure 5-4 now considers the entire area to be a 'cold' spot with GiZScores of less than -2.0.

Results such as these highlight the need for spatial researchers to examine exactly how the MAUP will affect their spatial research. The examples in Figures 5-1 to 5-4 do stress opposite ends of a distance threshold continuum, the lowest distance threshold compared to much higher distance threshold, though the results of this change are significant in several regards. These figures call attention to how drastic the scale effect of the MAUP can be in terms of how readily assemblages of spatial features can be subjectively illustrated and measured. Whether or not a group of features are considered "clusters" continues to depend on the scale with which the features are viewed. It is therefore imperative that any spatial research having to deal with cluster analysis consider employing multiple scales to analyze spatial data. Such considerations might promote awareness of the widespread effects of the MAUP. The rest of this section details some of the most noteworthy effects of the Cluster and Outlier analysis through analyses at

multiple scales and with different datasets in order to illustrate the resulting aggregation

effects.

Figure 5-3: Smaller-scale close up of Hotspot analysis map of Parcels3 with a 750 meter distance threshold.
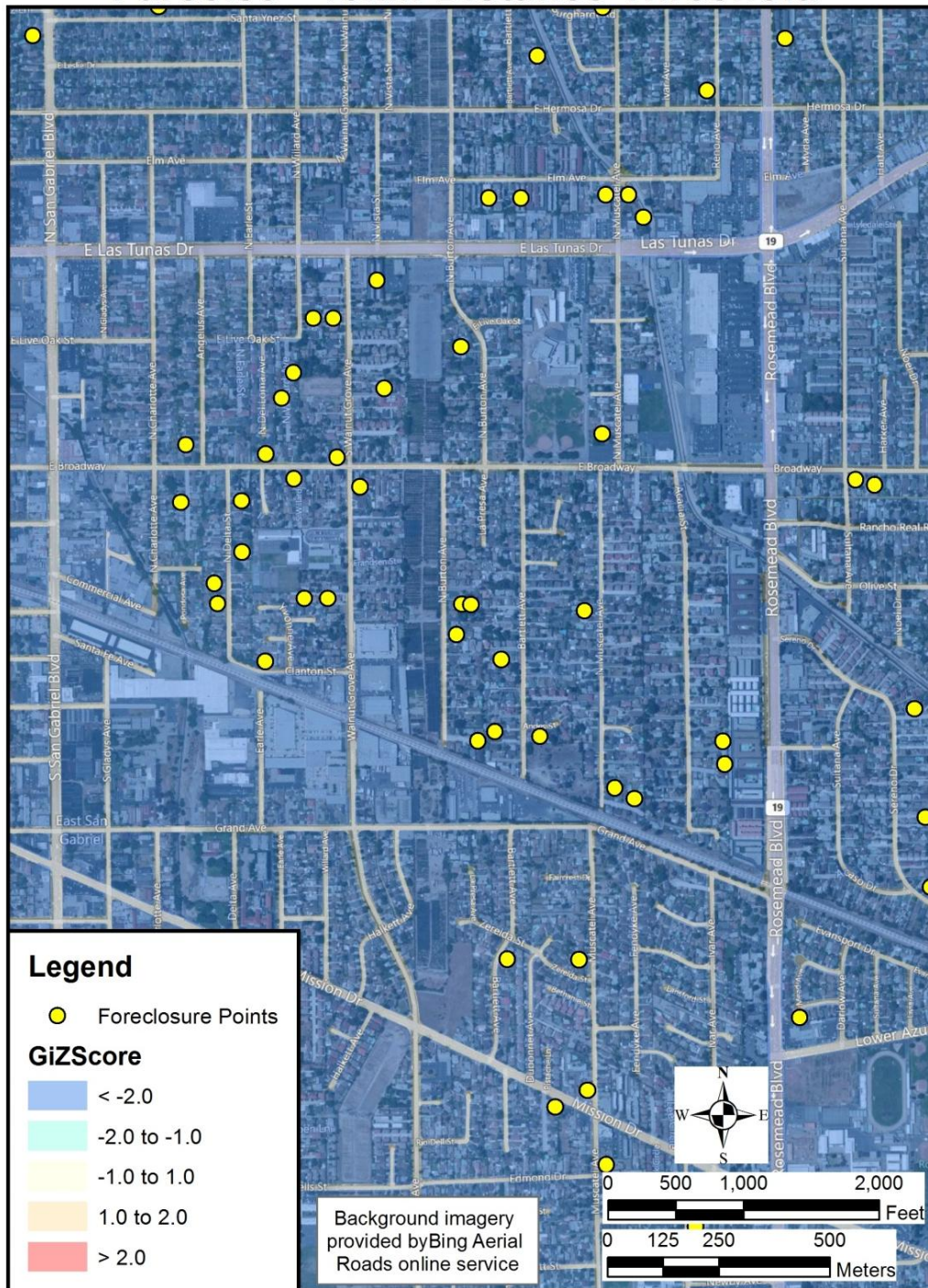
Figure 5-4: Smaller-scale close up of Hotspot analysis map of Parcels3 with a 10 kilometer distance threshold.

**5.2 Interpretation of Cluster and Outlier Analysis Results (Anselin Local Moran's I)**

The previous section dealt with what happens when the distance threshold changes in a given area in a hotspot analysis, whereas this section describes what happens when the distance threshold remains the same while the datasets vary. Such variation will highlight how scale/aggregation and zoning effects alter the results of a Cluster and Outlier analysis. When comparing Figures 5-5 and 5-6 or Figures 5-5 and 5-7 to each other, the shift depicted by the variation in hot, cold, or dispersed is indicative of how much zoning affects the results. Cluster and outlier analysis looks at how features relate to other neighborhood features. There are two blue cells above the legend in Figure 5-5. These features are dark blue not just because there is a lack of foreclosures within their boundaries, but also because surrounding cells containing more than one instance of foreclosure. When compared to Figures 5-6 and 5-7 their presence is completely lost when this information was aggregated into shifted zonal units.

Figures 5-8, 5-9, and 5-10 illustrate a similar point. Each of the figures, when compared to one another, gives different results. Due to the increasing aggregation of the Parcels features where tracts of land were merged into one another, some features are often only separated by a single street. A strong example of this can be viewed in the way that certain areas have several 'clusters' of foreclosures (depending on scale), though the Parcels dataset does not distinguish naturally occurring neighborhoods, merely (non)adjacent boundaries. The question of which scale and aggregation method work best is discussed in the next section.

# Cluster & Outlier Analysis (Close View)
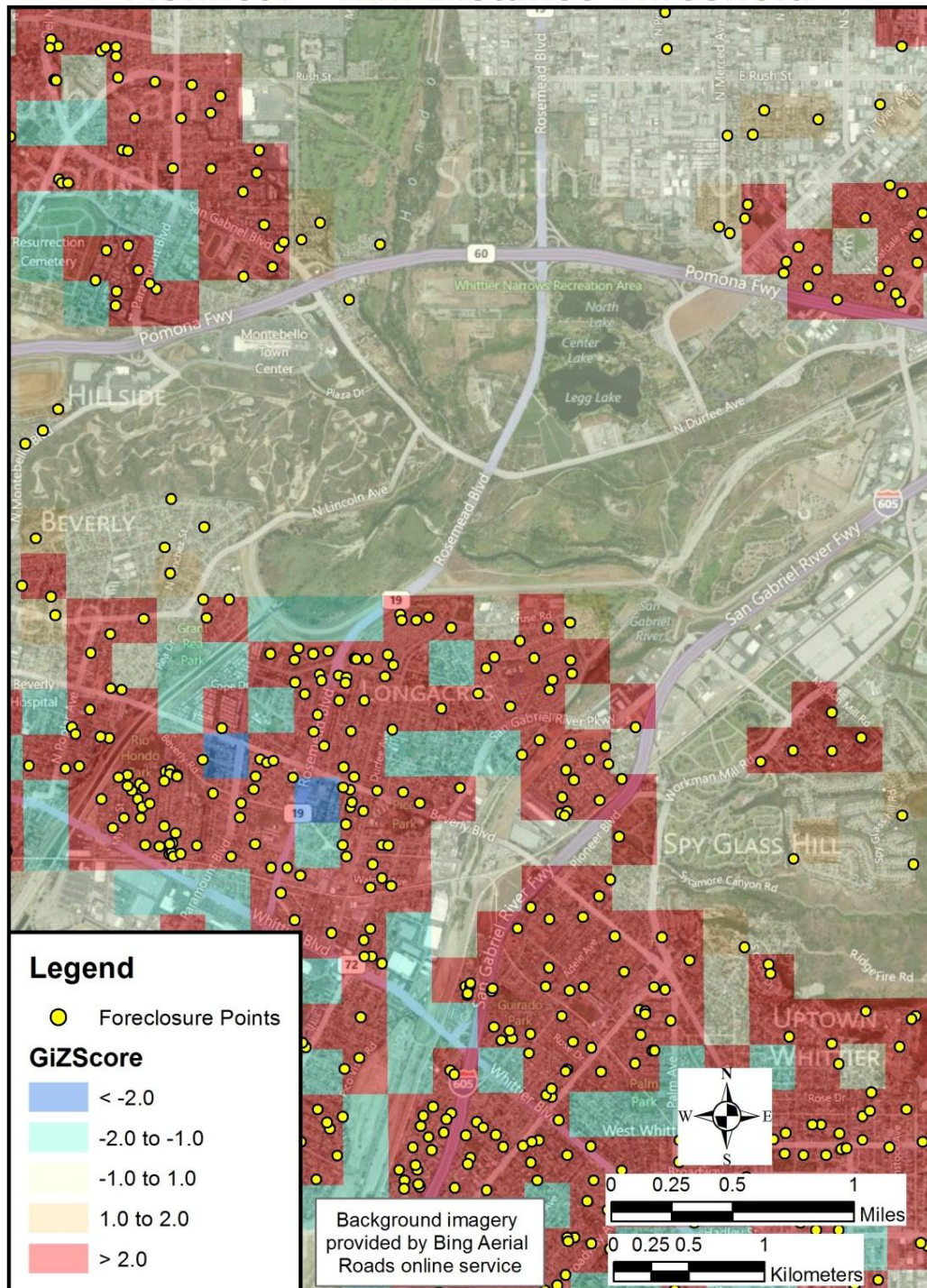## Fishnet1 - 1km Distance Threshold



Figure 5-5:  Cluster and outlier analysis map of Fishnet1 with a 1 kilometer distance threshold.

# Cluster & Outlier Analysis (Close View)
# Fishnet2 - 1km Distance Threshold



## Legend

○ Foreclosure Points

**GiZScore**

| | |
|---|---|
| | < -2.0 |
| | -2.0 to -1.0 |
| | -1.0 to 1.0 |
| | 1.0 to 2.0 |
| | > 2.0 |

Background imagery provided by Bing Aerial Roads online service
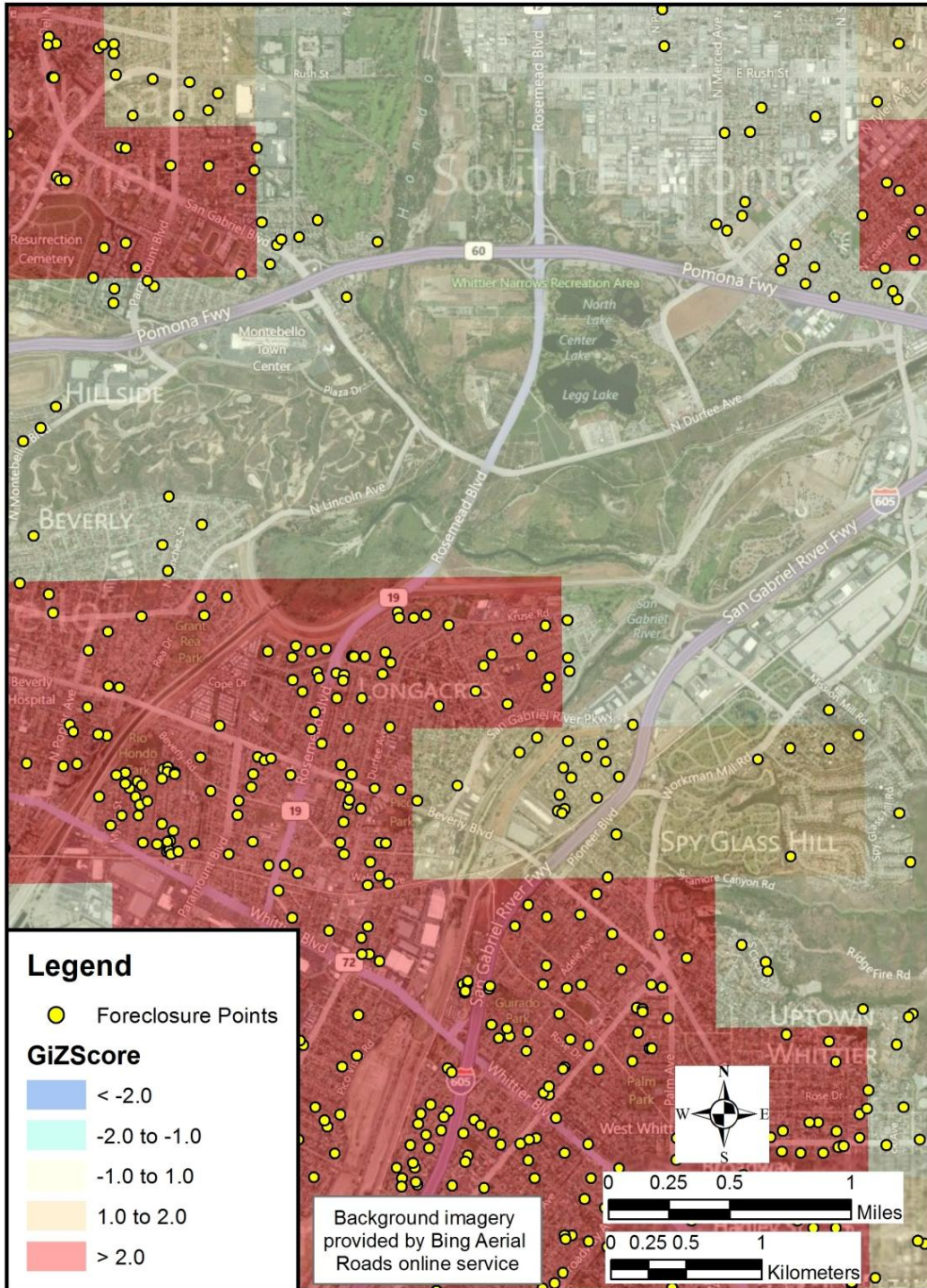
Figure 5-6: Cluster and outlier analysis map of Fishnet2 with a 1 kilometer distance threshold.

# Cluster & Outlier Analysis (Close View)
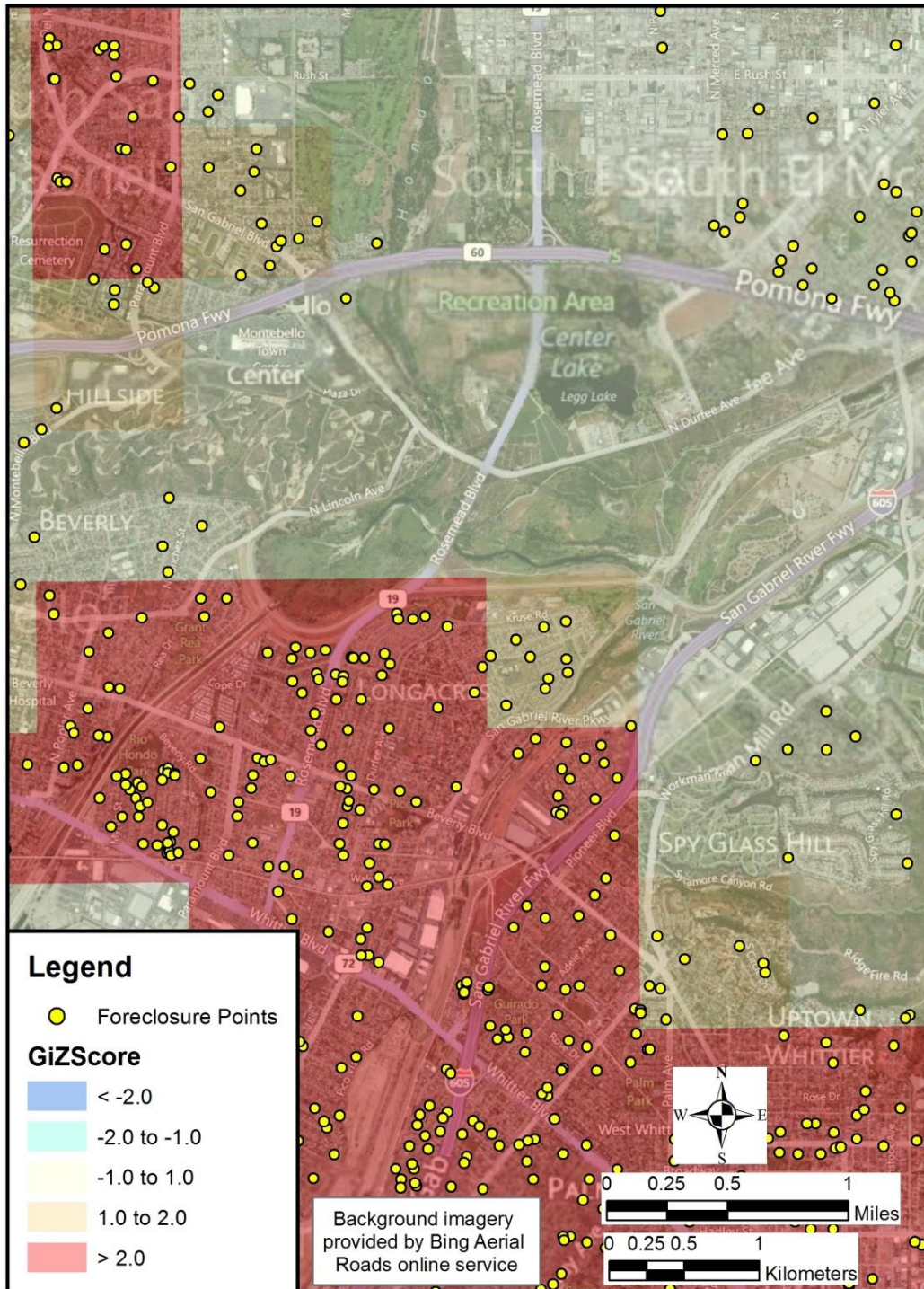# Fishnet3 - 1km Distance Threshold



Figure 5-7: Cluster and outlier analysis map of Fishnet3 with a 1 kilometer distance threshold.

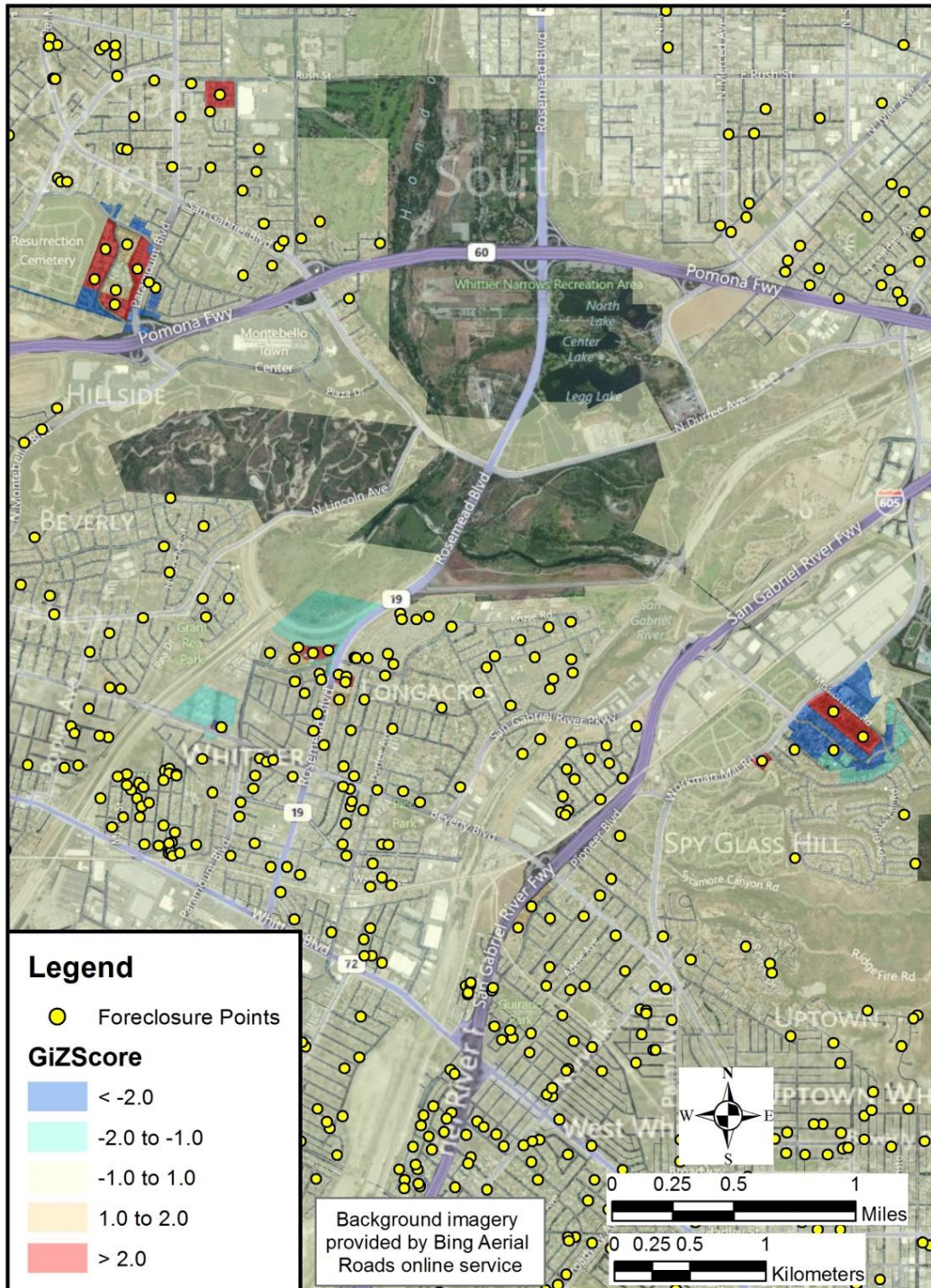# Cluster & Outlier Analysis (Close View)
# Parcels - 1km Distance Threshold



Figure 5-8: Cluster and outlier analysis map of Parcels with a 1 kilometer distance threshold.

# Cluster & Outlier Analysis (Close View)
## Parcels2 - 1km Distance Threshold



**Legend**

⬤ Foreclosure Points

**GiZScore**

| | |
|---|---|
| | < -2.0 |
| | -2.0 to -1.0 |
| | -1.0 to 1.0 |
| | 1.0 to 2.0 |
| | > 2.0 |

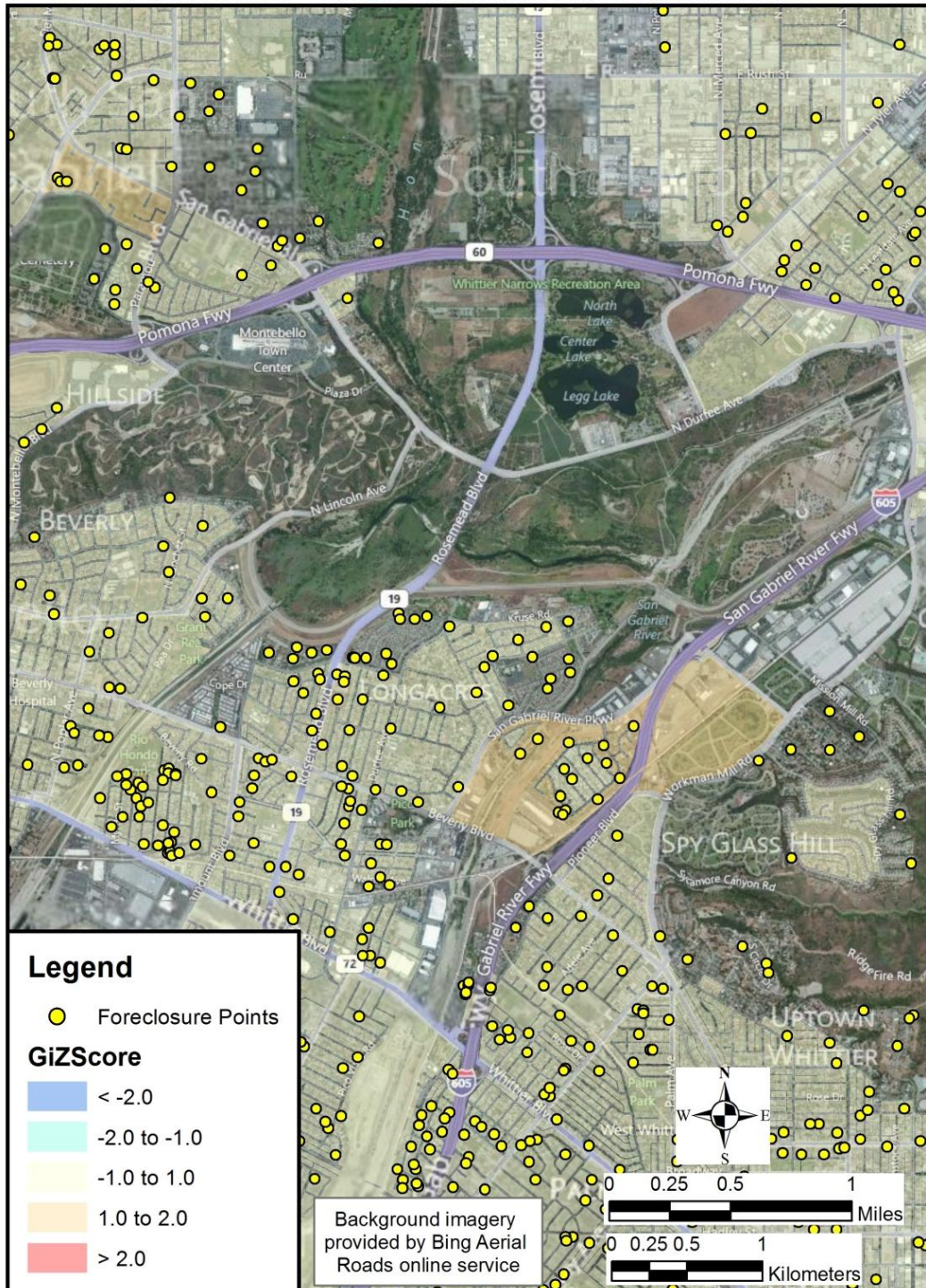Background imagery provided by Bing Aerial Roads online service

Figure 5-9: Cluster and outlier analysis map of Parcels2 with a 1 kilometer distance threshold.
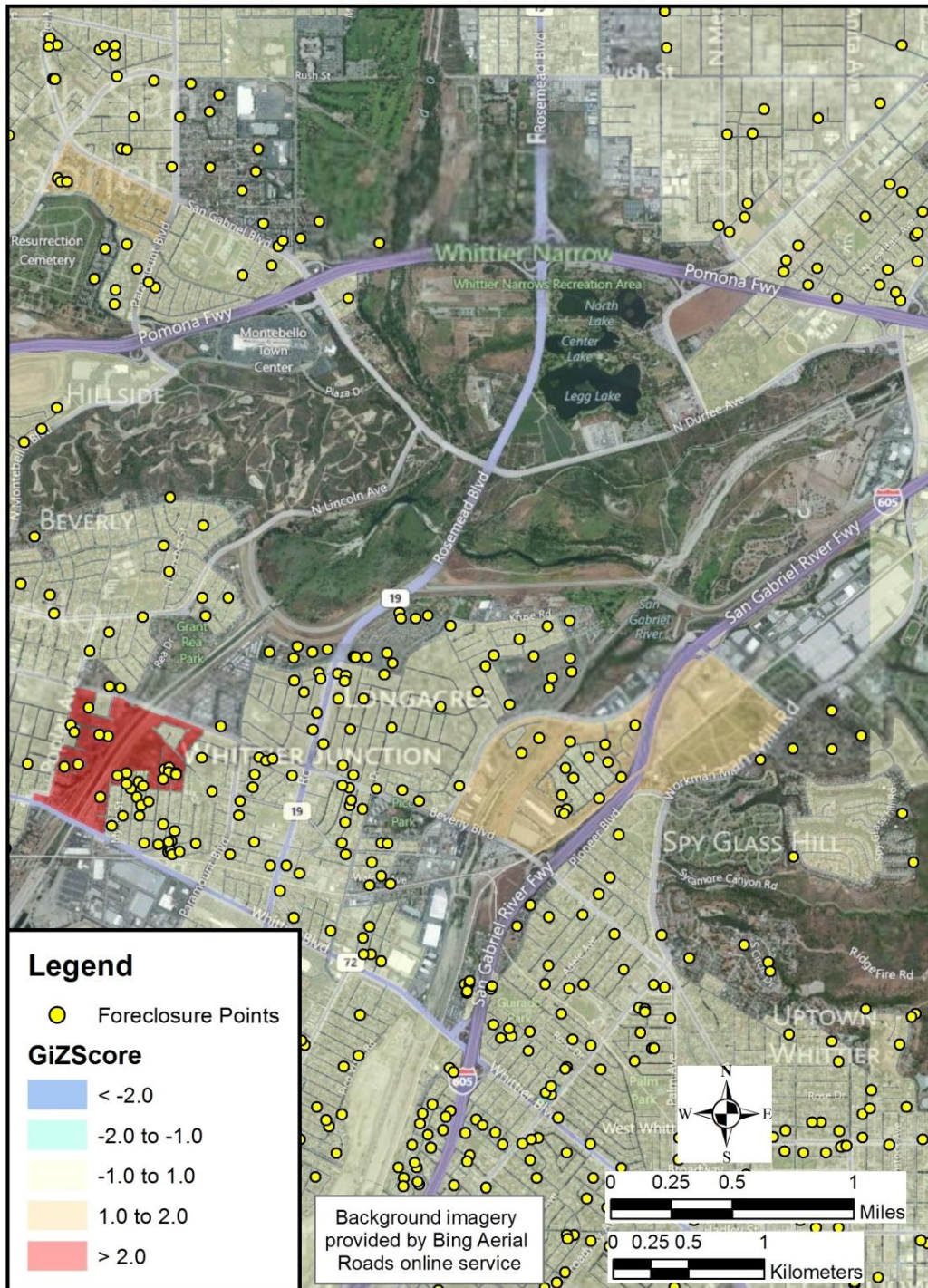
Figure 5-10: Cluster and outlier analysis map of Parcels3 with a 1 kilometer distance threshold.

**5.3 The Broader Results**

The overall pattern of the results indicates strong clustering in the foreclosure

points and the various aggregations of those points.  The tools which analyze the global

patterns of clustering and dispersion clearly indicate a prevalence of clustering across

multiple datasets and distance thresholds.  When reviewing the mapped results of the

hotspot analysis and the cluster and outlier analysis, results show the strongest instances

of variation with the datasets.  Areas with high distance thresholds are shown to be

clustered while the same areas indicate disorder when lower distance thresholds are

chosen.  This type of variation exists across all mapped datasets where any single input is

altered.  There is very little homogeneity between and amongst the different mapped

outputs.

The implications of these varying outputs are far reaching.  They not only confirm

the existence of the MAUP but also how difficult it can be to determine the most

appropriate approach spatial research would need to take to account for its effects.  These

outputs highlight the importance of running various types of analyses to make certain that

results accurately reflect the true physical topic of interest, and not skewed results that

aggregation and zoning effects will create.

As a hypothetical example, a city planner wants to know what areas of a city need

more street lights for its citizens. The planner would need to delineate how he or she

characterizes the areal units in this type of study.  Are they constructed of parcel units,

regular grids, commonly-accepted neighborhood delineations, or some other artificially

constructed boundary?  Perhaps it is a combination of multiple types of boundaries in

order to find the least luminous areas.  Does the planner also account for population density so that more lights are needed in denser areas?  It would also need to be determined how many lights a given areal unit would need in order to be considered sufficiently luminous.  This hypothetical example is riddled with questions that need to be answered by generating multiple outputs to thoroughly answer the questions the analysis puts forward.

## 5.4 Future Avenues of Research

The scope of this project is based upon the exemplification of the effects of the MAUP using cluster analysis methods on instances of real foreclosures in Los Angeles County, between 2006 and 2008, and various aggregations of these events.  There are a number of avenues of exploration remaining for future appropriate inquiry into this area of research.  This project researched only instances of foreclosures and their spatially joined counts as the primary variable studied.  In the future, it would be interesting to examine the regression and autocorrelation results of the foreclosure data with demographic data, such as household income, residential age groups, and crime rates.  In addition to such analyses, yet more research possibilities include determining the 'true' aggregation effects when factoring in non-Euclidian distance measures.  An example of this approach would be to include two foreclosures happening within 500 meters from one another, but separated geographically by elevation.  Another possible avenue of research may be in the replication of cluster analysis results with point features from highly different fields of spatial research, such as locations of infected flora or locations

of street lights signifying 'clusters' of illumination and thus light pollution. Although these additional avenues of research could be used to further validate the analyses of the MAUP and shed more light on the best practices, they do not necessarily encompass the scope of a typical thesis project.

## 5.5 Answering the Theoretical Questions

Chapter One introduced the modifiable areal unit problem, or MAUP, and the scale and aggregation problems which comprise it. Subsequent sections gave information on just how much of an impact scale and aggregation have on analysis results. It was determined from the results of this research project that issues of scale and zoning affect conceptualization as well as the results of such analyses. The scale of what is utilized – large/small, many/few, coarse/fine, and high/low distance thresholds – impacts how a project is formulated and executed, and in turn the results.

To answer the first thesis question of this project, the effects of the MAUP on the visual results are clearly shown in the maps and stress the importance of looking at wide ranges of areal scales and zonal configurations in order to best fit avenues of research. It is challenging to navigate the visual differences in the maps of areas of Los Angeles County that result from the cluster analysis tools and their varying datasets and distance thresholds.

The second thesis question, regarding best scale or zonal unit of analysis, assumes that all spatial research can and should use a standard zonal scheme. This project looked at regularly gridded areas at different scales and zonal configurations (Fishnet1, Fishnet2,

and Fishnet3) as well as more realistic structures that represent actual urban geographies (Parcels, Parcels2, Parcels3). The unit selected should be comparable to the size of the phenomenon being studied. Researchers have to use intuition and educated guessing to know what this is or use a variety of scales and see which one gives results that are close to what should be expected. The questions that each spatial researcher attempt to answer require multifaceted approaches to deal with the scale effect and the zoning effect.

More research should be conducted to highlight known methods of MAUP mitigation. One of the strongest conclusions that can be drawn from this research is that most forms of spatial analysis must be sure to at least acknowledge the assorted impacts of the MAUP by varying analyses at different scales relevant to the given study area, aggregation units, and perspectives in order to fully comprehend all the possible consequences of the MAUP on their research.

The final question this research attempted to answer was how the MAUP affects the way in which a methodology should be constructed. Any and all measures should be taken in order to account for the MAUP in research projects. Steps should always be taken to mitigate the effects of the MAUP in spatial research (Mu and Wang, 2008). Certain approaches, such as the 'moving window' approach, can be used to help mitigate some aspects of the MAUP in spatial research (Jelinski and Wu, 1996; Ratcliffe and McCullagh, 1999). Acknowledging the MAUP in published research can go a long way towards cultivating the thought processes needed to ensure the future development of sound methodologies which account for the inherent biases of preselected boundaries, scales, extents, and the limitations of aggregated data in spatial research (Hipp, 2007).

The surest way of imparting the impacts of the MAUP on those performing spatial research is through rigorous education and training. The development of GIS technology coupled with the ever-increasing advances in computer and network capabilities is generating spatial data at an unprecedented rate. One need only look at the various spatial data repositories on the internet to recognize the variety of potential avenues of analysis available to everyone.

Persistent effects generated by the MAUP will continue to plague spatial research until large scale efforts are taken to understand, quantify, and evaluate all facets of the MAUP. An even greater effort is required to educate those who would perform spatial research on how to anticipate the MAUP and mitigate its effects, if possible. It must be proposed that the first step in such educational efforts include mention of the MAUP in more research and literature than is currently present. Perhaps such educational efforts will include the introduction of the background and prevalence of the MAUP. When enough researchers and students are made aware of the MAUP and its effects become known to a wider audience it will then permeate guidebooks, textbooks, and the technical literature which will influence generations of GIS professionals. Such permeation will saturate the GIS community which will enable understanding of the MAUP and hopefully generate a level of enthusiasm needed to help mitigate it on an institutional level.

# References

Amrhein, C.G. 1993. Searching for the Elusive Aggregation Effect: Evidence From Statistical Simulations. *Environment and Planning A*, Vol. 27, pp. 105-119.

Amrhein, C.G., and H. Reynolds. 1996. Using Spatial Statistics to Assess Aggregation Effects. Geographical Systems, Vol. 2, pp.83-101.

Amrhein, C.G., and H. Reynolds. 1997. Using the Getis Statistic to Assess Aggregation Effects in Metropolitan Toronto Census Data. The Canadian Geographer, Vol 31, No. 2, pp. 137-149.

Bhati, Avinash Singh. 2005. Robust Spatial Analysis of Rare Crimes: An Information-Theoretic Approach. *Sociological Methodology*, Vol. 35, pp. 239-301.

Ceccato, Vania, Robert Haining, and Paola Signoretta. 2002. Exploring Offence Statistics in Stockholm City Using Spatial Analysis Tools. *Annals of the Association of American Geographers*, Vol. 92, No. 1, pp. 29-51.

Chainey, Spencer, L. Thompson, and S. Uhlig. 2008. The Utility of Hotspot Mapping for Predicting Spatial Patterns of Crime. *Security Journal*, Vol. 21. Pp 4-28.

Cressie, Noel and Linda B. Collins. 2001. Analysis of Spatial Point Patterns Using Bundles of Product Density LISA Functions. *Journal of Agricultural, Biological, and Environmental Statistics*, Vol. 6, No. 1. 118-135.

Dale, Mark R.T., P. Dixon, M.-J. Fortin, P. Legendre, D.E. Myers, and M.S. Rosenberg. 2002. Conceptual and Mathematical Relationships among Methods for Spatial Analysis. *Ecography*, Vol. 25, No. 5, pp. 558-577.

Dark, Shawna J. and Danielle Bram. 2007. The Modifiable Areal Unit Problem (MAUP) in Physical Geography. *Progress in Physical Geography*, Vol. 35, No. 5, pp. 471-479.

83

Diggle, P.J., P.J. Ribeiro Jr., and O.F. Christensen. 2003. An Introduction to Model-based Geostatistics. Chapter 2: Møller, J. (ed.) *Spatial statistics and computational methods.* Springer Verlag, 2003.

Fotheringham, A.S. 1989. Scale-Independent Spatial Analysis. In *Accuracy of Spatial Databases*, edited by M.F. Goodchild and S. Gopal, London: Taylor and Francis. Pp 221-228.

Fotheringham, A.S. and Wong, D.W.S., 1991, The Modifiable Areal Unit Problem in Multivariate Statistical Analysis. *Environment and Planning A*, Vol. 23, 1025-1044.

Gatrell, Anthony C., T.C. Bailey, P.J. Diggle, and B.S. Rowlingson. 1996. Spatial Point Pattern Analysis and its Application in Geographic Epidemiology. *Transactions of the Institute of British Geographers. New Series*, Vol. 21, No. 1, pp. 256-274.

Getis, Arthur and Janet Franklin. 1987. Second-Order Neighborhood Analysis of Mapped Point Patterns. *Ecology*, Vol. 68, No. 3, pp. 1 473-477.

Getis, Arthur and K. Ord. 1992. The Analysis of Spatial Association by Use of Distance Statistics. Geographical Analysis, Vol. 24, pp. 189-206.

Gotway, Carol A. and Linda J. Young. 2002. Combining Incompatible Spatial Data. *Journal of the American Statistical Association*, Vol. 97, No. 458, pp. 632-648.

Haase, Peter. 1995. Spatial Pattern Analysis in Ecology Based on Ripley's K-Function: Introduction and Methods of Edge Correction. *Journal of Vegetation Science*, Vol. 6, No. 4. Pp. 575-582.

Hay, G.J., D.J. Marceau, P. Dube, and A. Bouchard. 2001. A Multiscale Framework for Landscape Analysis: Object-Specific Analysis and Upscaling. *Landscape Ecology*, Vol. 16, pp. 471-490.

Hayward, Peter and Jason Parent. 2009. Modeling the Influence of the Modifiable Areal Unit Problem (MAUP) on Poverty in Pennsylvania. *The Pennsylvania Geographer*, Vol. 47, No. 1. Pp 120-135.

Hipp, John R. 2007. Block, Tract, and Levels of Aggregation: Neighborhood Structure and Crime and Disorder as a Case in Point. *American Sociological Review*, Vol. 72, No. 5, pp. 659-680.

Jelinski, Dennis E. and Jianguo Wu. 1996. The Modifiable Areal Unit Problem and Implications for Landscape Ecology. *Landscape Ecology*, Vol. 11, No. 3, pp. 129-140.

Larsen, Frank W. and Carsten Rahbek. 2003. Influence of Scale on Conservation Priority Setting – a Test on African Mammals. *Biodiversity and Conservation*, Vol. 12, pp. 599-614.

Lentz, Jennifer A., Jason K. Blackburn, and Andrew J. Curtis. 2011. Evaluating Patterns of a White-Band Disease (WBD) Outbreak in Acropora palmate Using Spatial Analysis: A Comparison of Transect and Colony Clustering. *PloS ONE*, Vol. 6, No. 7, pp. 1-10.

MacEachren, Alan M. 1982. Choropleth Map Accuracy: Characteristics of the Data. *Technical Papers of ACSM, Denver*, pp. 512-521.

Mitchell, Andy. 2009. The ESRI Guide to GIS Analysis, Volume 2: Spatial Measurements and Statistics. *ESRI Press*. Redlands, California.

Mueller-Warrant, George W., G.W. Whittaker, and W.C. Young III. 2008. GIS Analysis of Spatial Clustering and Temporal Change in Weeds of Grass Seed Crops. *Weed Science*. Vol. 56, No. 5, pp. 647-669.

Mu, Lan and Fahui Wang. 2008. A Scale-Space Clustering Method: Mitigating the Effect of Scale in the Analysis of Zone-Based Data. *Annals of the Association of American Geographers*, Vol. 98, No. 1, pp. 85-101.

Nakaya, Tomoki. 2000. An Information Statistical Approach to the Modifiable Areal Unit Problem in Incidence Rate Maps. *Environment and Planning*, Vol. 32, Pp. 91-109.

Openshaw, Stan. 1984. *The Modifiable Areal Unit Problem*. CATMOG #38. Norwich: Geo Books.

Openshaw, S. and Taylor, P. J. 1979. A Million or So Correlation Coefficients: Three Experiments on the Modifiable Areal Unit Problem. In Wrigley, N., ed., *Statistical methods in the spatial sciences*. London: Pion, 127-44.

Pawitan, Gandhi and David G. Steel. 2009. Exploring the MAUP From a Spatial Perspective. *Center for Statistical and Survey Methodology, University of Wollongong, Working Paper* 20-09, p. 1-28.

Pendergrast, S., N. Wood, J.H. Lawton, and B.C. Eversham. 1993. Correcting for Variation in Recording Effort in Analyses of Diversity Hotspots. *Biodiversity Letters*, Vol. 1, No. 2, pp. 39-53.

Perry, George L.W., B.P. Miller, and N.J. Enright. 2006. A Comparison of Methods for the Statistical Analysis of Spatial Point Patterns in Plant Ecology. *Ecology*, Vol. 187, No. 1, pp. 59-82.

Rahbek, Carsten and Gary R. Graves. 2000. Detection of Macro-Ecological Patterns in South American Hummingbirds is Affected by Spatial Scale. *Proceedings: Biological Sciences*, Vol. 267, No. 1459, pp. 2259-2265.

Ratcliffe, J.H. and M.J. McCullagh. 1999. Hotbeds of Crime and the Search for Spatial Accuracy. *Journal of Geographical Systems*, Vol. 1, pp. 385-398.

Ripley, B.D. 1977. Modelling Spatial Patterns. Journal of the Royal Statistical Society. Series B (Methodological). Vol. 39, No. 2, pp. 172-212.

Rushton, G. 1998. Improving the geographic basis of health surveillance using GIS. In *GIS and Health*. Taylor and Francis, London, pp. 63-79.

Rushton G. and P. Lolonis. 1996. Exploratory Spatial Analysis of Birth Defect Rates in an Urban Population. *Statistics in Medicine*, Vol. 7, pp. 717-726.

Shriner, Susan A., K.R. Wilson, and C.H. Flather. 2006. Reserve Networks Based on Richness Hotspots and Representation Vary with Scale. *Ecological Applications*, Vol. 16, No. 5, pp. 1660-1671.

Tagashira, Naoto and Atsuyuki Okabe. 2002. The Modifiable Areal Unit Problem in a Regression Model Whose Independent Variable Is a Distance from a Predetermined Point. *Geographical Analysis*, Vol. 34, No. 1, pp. 1-20.

Williamson, D., McLafferty, S., McGuire, P., Ross, T., Mollenkopf, J., Goldsmith, V. and Quinn, S. 2001. Tools in the Spatial Analysis of Crime. In *Hirshfield, A. and Bowers, K. (eds) Mapping and Analyzing Crime Data: Lessons from Research and Practice. London: Taylor and Francis*, pp. 187-202.