

# A Comparison of Methods for Spatial Interpolation across Different Spatial Scales

Idris Mohammed Jega<sup>1</sup>, Alexis J. Comber<sup>2</sup>, Nicholas J. Tate<sup>3</sup>

<sup>1</sup>Strategic Space Applications, National Space Research and Development Agency, Abuja, Nigeria

<sup>2</sup>School of Geography, University of Leeds, Leeds LS2 9JT, UK

<sup>3</sup>Department of Geography, University of Leicester, Leicester LE1 7RH, UK

## Abstract

*Spatially distributed estimates of population provide commonly used demand surfaces in support of spatial planning. In many countries, spatially detailed population estimates in small areas are not available. For such cases a number of interpolation methods have been proposed to redistribute summary population totals over small areas to estimate locally nuanced demand surfaces. Population allocations to small areas are commonly validated by comparing the estimates with some known values for those areas. This paper explores different interpolation methods applied at different spatial scales in locations where the validation of estimated surfaces is possible in order to suggest appropriate interpolation parameters for locations where it is not. The results show binary dasymetric mapping applied at medium scales provide the best estimates of population, among the methods, areal weighting the worst at all scales and pycnophylactic interpolation shows significant improvement on areal weighting at all scales. This paper provides a comprehensive evaluation of these techniques, using different scales of input data and residual mappings to compare and evaluate the spatial distribution of errors in the estimated surfaces. The application of such methods for estimating spatially distributed demand population values in different types of spatial data analysis and in locations where validation data do not exist are discussed.*

**Keywords** areal interpolation, dasymetric mapping, areal weighting, pycnophylactic interpolation, population data.

## I. INTRODUCTION

Population estimates for small areas are important for many types of spatial data analysis. They are especially important for accessibility studies and facility location-allocation analyses, both of which are commonly used to support spatial planning and policy development. Population censuses provide a reliable record of socioeconomic characteristics and the spatial distribution of residential population [1] and thereby support geodemographic analyses [2]. In the U.K. census, population counts are collected for each household and published as aggregate counts and statistics for fixed pre-defined spatial units with Output Areas (OA) being the most detailed. The OA is similar to a U.S. census block. The OA was designed to be as homogenous as possible and to

have a similar population size [3,4]. The target size of an OA is 125 households or approximately 300 people [3]. The main reason for aggregating population census counts in this way is to maintain confidentiality and respondent anonymity. In some countries census data are spatially aggregated only to very coarse summaries that limit their use in further spatial analysis. For example, in Nigeria, simple population totals are provided for each state and local government areas (LGAs), with the LGA providing the most spatially detailed information. An LGA is similar to the size of a county or Unitary Authority (UA) district in the U.K. This level of aggregation makes many types of spatial data analysis difficult because more detailed population estimates are often required than those provided [5].

## II. BACKGROUND

Areal interpolation is the process of transforming values of interest from *source zones* to provide estimates over a set of *target zones* with unknown values [6]. A number of areal interpolation techniques have been developed and their performance has been found to relate to specific characteristics of the input data including its errors, extent and spatial properties [7,8], as well as the characteristics of any ancillary data used, for example, to constrain the disaggregation [1].

One of the simplest areal interpolation techniques is areal weighting. In this total data volumes are maintained under the assumption that population is uniformly distributed within the source zones [6]. In reality, population distributions are not uniform within source zones and assigning the same population density to every location may not represent the actual population distribution because of the presence of unpopulated areas (water bodies, parks, industrial areas, etc.). Point-based areal interpolation methods [9] have been used to overcome some of the errors associated with the assumption of uniform densities within source zones. These methods assign census zone populations to the centroid of each source zone, and then population counts are estimated by summing all points within the target zone. The major shortcoming of this method is that the polygon centroid is used to represent the total population within the polygon. When the source and target zones are spatially intersected, the total population is completely allocated (or not) to the target zone, depending on location of the centroid

[10]. This can cause errors when the populations allocated in this way are used as demand surfaces for measuring access to service facilities [11]. Tobler [12] proposed pycnophylactic interpolation as a technique to overcome this shortcoming. These generate spatially varying but smooth surfaces, whilst preserving the total data volumes and assign a non-zero population density value to target zone. In reality, the target zones may have sudden changes in population density that coincide uninhabited areas. Thus approaches that make use of ancillary data to constrain areas within source zones have been suggested [13,14,15] and ancillary data on urban extent has been commonly used.

Remotely sensed data such as aerial photographs have been used to map urban extent since the 1950s and estimate populations [16]. Lo [17] describes three main approaches used to visually interpret aerial photographs for population mapping: counting individual dwelling units, extracting the extent of urban settlement and measuring areas of different land use. Digital images and statistical classification of broad land use types are now common [18,19], and land cover derived in this way has been used as ancillary data in spatial interpolation [1,15,20].

The dasymetric mapping approach is an areal interpolation technique that incorporates ancillary data sources as control variables in order to identify zones having different population densities [20,21]. It constrains the disaggregation of population values from source zones to specific target zones, which can be weighted by for example expected residential density [13,14,15,20,22,23]. Binary dasymetric approach [20] divides source zones into populated and unpopulated areas and allocates population only to the populated areas. Su et al. [24] extended this idea by further dividing the populated area into multiple classes using transportation layers, topography and land use zoning. A 3-class dasymetric model has been proposed [14,25], but has not been shown to provide any additional benefit to binary dasymetric approaches. Recent research has improved areal interpolation approaches by applying simple proportions as well as various forms of regression analysis [26], quantile regression [21] and through improved ancillary data such as LiDAR [27], open access vector map data [1] and household survey data [5].

In many interpolation studies population totals are redistributed from an initial area, the source zone (e.g. MSOA in the U.K.) to smaller target zones such as OAs and the results are compared with known population counts at the lower level in order to validate the method. Additionally, much previous research has used multispectral imagery mainly of 30m spatial resolution to redistribute aggregate census data to a lower level census unit as the target

zones for which true populations are known [15,17,24,25]. Langford [1] draws attention to the implications of this practice: first, the performance of the most spatially detailed census data are not often measured because they are reserved for testing the performance of the interpolation methods; and second, it is difficult to evaluate the performance of target zones smaller than the lowest level census spatial unit because their true values are not known. He demonstrates the possibility of using unit postcodes (UPCs) in the UK as the target zones with an acceptable precision. The UPCs are smaller than the finest census zone division, the OA in the U.K. The population totals of the UPCs are not reported in the U.K. hierarchy of census units but are known and available at the Office of National Statistics (ONS) U.K.

This study evaluated areal weighting, dasymetric mapping and pycnophylactic interpolation applied across different spatial scales and using ancillary land cover data classified from satellite imagery of differing spatial resolutions. Different interpolation methods and input parameters were applied to a U.K. case study to determine how well the populations reported in census small areas were estimated, and thus how well population values generated in this way could be used for accessibility studies, location-allocation analyses etc. It sought to address two specific research questions:

- The relationship between estimated populations from different interpolations and the known census counts?
- Which is the most appropriate interpolation method to apply in the absence of a universally accepted methodology in estimating population surfaces?

### III. MATERIALS AND METHODS

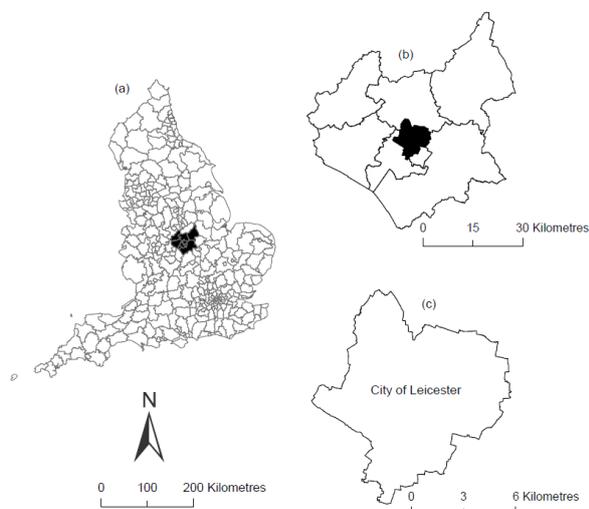
#### A. Study area

The study area was the city of Leicester in the UK, a location where the actual population distribution is known and where interpolation model output validation was possible. Leicester covers an area of about 73 km<sup>2</sup>. The population of Leicester has increased between 1990 and 2011 as shown in Table I. Figure 1 shows the location and extent of Leicester in the county of Leicestershire, in England. This location was chosen because of its proximity and the authors' knowledge of the area. The 2001 population data were used in this study.

**Table I Percentage Change in Population for Leicester from 1951 to 2011**

Census Year	Population
1951	285200
1961	288100
1971	284200
1981	280300

1991	272133
2001	279921
2011	329839



**Fig 1: The map of (a) England Showing Location of Leicestershire County; (b) Leicestershire County with Location of Leicester UA; (c) Leicester UA. The Digital Boundaries are © Crown Copyright and/or database right 2013. An Ordnance Survey/EDINA Supplied Service**

**B. Data**

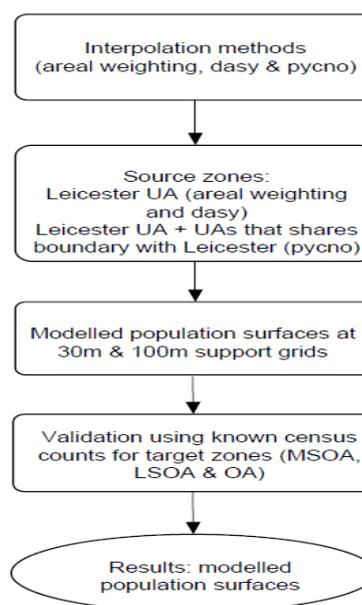
The aim of this research is to develop a novel and comprehensive analysis of the operation of three classic spatial interpolation approaches and how they interact with different target zone sizes, support grids and different scales of ancillary data. The source zone was the city of Leicester, a unitary authority administrative area and the target zones evaluated were, in order of increasingly granularity, MSOAs, LSOAs and OAs. Satellite imagery covering the study area was acquired to generate land cover data and support the dasymetric approaches. Medium resolution satellite imagery and fine resolution (25cm aerial photography) were used to generate ancillary data for the dasymetric analyses. Table II summarises the data used in this study.

**Table II Data for the city of Leicester**

Data	Format	Date	Source
Landsat7 (ETM) 30m spatial resolution	Image	16 April 2003	United States Geological Survey (USGS) website ( <a href="http://www.usgs.gov/">http://www.usgs.gov/</a> )
Ortho-rectified aerial photograph 25cm spatial resolution	Image	22 May 2010	Ordnance Survey, U.K. © Crown copyright and/or database right 2013. All rights reserved.
Census data with boundaries of OAs, LSOAs and MSOAs	Shapefile	2001 Census	Census Area Statistics on the Web (casweb) ( <a href="http://casweb.mimas.ac.uk/2001/start.cfm">http://casweb.mimas.ac.uk/2001/start.cfm</a> ).

**C. Analysis**

An overview of the analysis is shown in Figure 2. Areal weighting and dasymetric methods were applied to Leicester unitary authority as the source zone. The pycnophylactic interpolation method was applied to census totals for Leicester unitary authority together with all the surrounding unitary authorities (Harborough, Blaby, Charnwood and Oadby and Wigston) to generate an interpolated gridded population surface at resolutions of 100m and 30m which were then summed over MSOA, LSOA and OA target areas. This is because the pycnophylactic method cannot be applied to a single polygon such as the Leicester unitary authority. The estimated populations were then compared with the known census counts in each case, for validation.



**Fig 2: An Overview of the Method**

1) **Areal Weighting**

Areal weighting is based on the assumption that the true population is uniformly distributed within source zones [6]. It uses the size (area) of each target zone to proportionally allocate the population. It was implemented in six steps: (1) the area of the source zone was calculated; (2) the population density was calculated using Equation 1; (3) the source zone and target zones were spatially intersected; (4) the intersect areas were calculated; (5) a population value for each intersected zone was calculated from its area and the population density as in Equation 2; and (6) The interpolated population estimate for each target zone was calculated by summing all intersected areas within each target zone. A flowchart describing these steps is shown in Figure 3. The population density of the source zone is expressed mathematically as:

$$d_{sp} = \frac{P_s}{A_{sp}} \quad [1]$$

Where  $d_{sp}$  is the population density of the source zone,  $P_s$  is the total population of source zone  $s$  and  $A_{sp}$  is the area of source zone  $s$ .

The estimated population for the overlaid zones is expressed mathematically as:

$$\bar{P}_t = \sum_{s=1}^s A_{t,sp} d_{sp} \quad [2]$$

Where  $\bar{P}_t$  is the estimated populations of overlaid zone,  $t$ ;  $s$  is the number of source zones,  $A_{t,sp}$  is the area of intersection between overlaid zone  $t$  and source zone  $s$  and  $d_{sp}$  as defined in Equation 1.

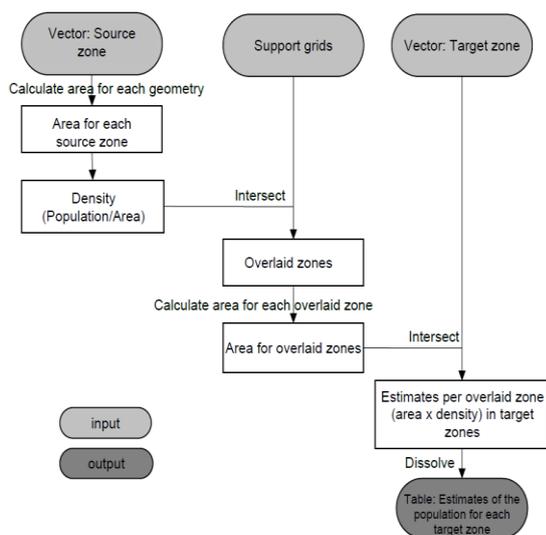


Fig 3: Implementation Steps for Areal Weighting (Vector Mode)

2) **Binary Dasymetric**

The binary dasymetric uses ancillary data to spatially constrain the disaggregation. It follows the same steps as above but the source zone and target zone areas are modified by removing non-populated areas from the analysis. The result is that the population density is calculated by dividing the population count of the source zone by the total size of all built-up areas within the source zone. In this case urban / non-urban areas were identified from Landsat7 (ETM) 30m spatial resolution data. This was classified into 3 classes Built-up, Water and Vegetation, with the first class forming the urban area and the last 2 the non-urban areas. The 25cm orthorectified aerial photography was used to resample image pixels to 3m and 10m without altering the projected coordinate system. Cubic convolution resampling was used to compute each output cell value because this method reduces blurring and produces a smoother output image than other commonly used method such as nearest neighbour or bilinear interpolation. The resampled images were classified to derive land cover data of the Leicester area at 10m and 3m spatial resolution. A supervised maximum likelihood classification identified the built-up areas and was repeated several times and the data with the highest classification accuracy of 87.89%, 83.20% and 82.03% for 30m, 10m and 3m resolution respectively were chosen. Accuracy of classification was assessed by comparing 256 randomly chosen pixels for which the land cover class was known. Figure 4 shows the classified Landsat7 (ETM) data and Figure 5 shows the binary mask derived from the classified image. The application of the binary dasymetric method is summarised in Figure 6.

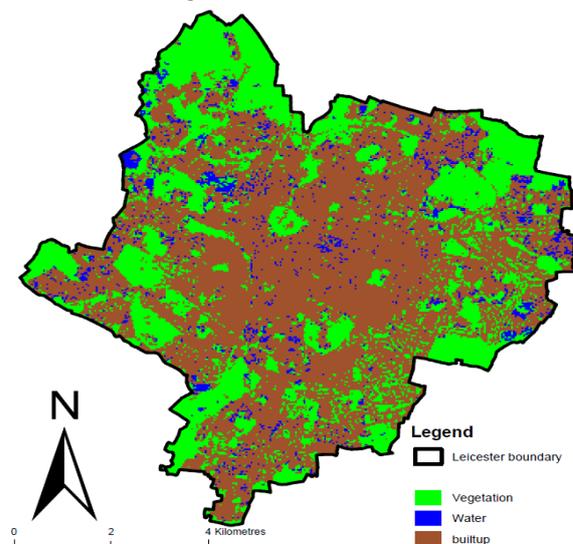
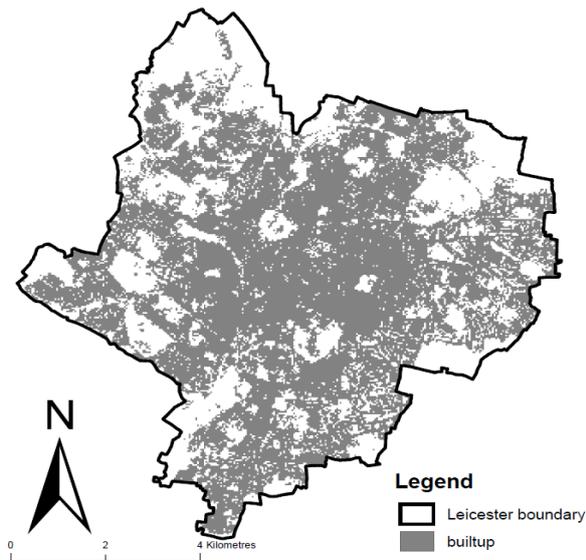


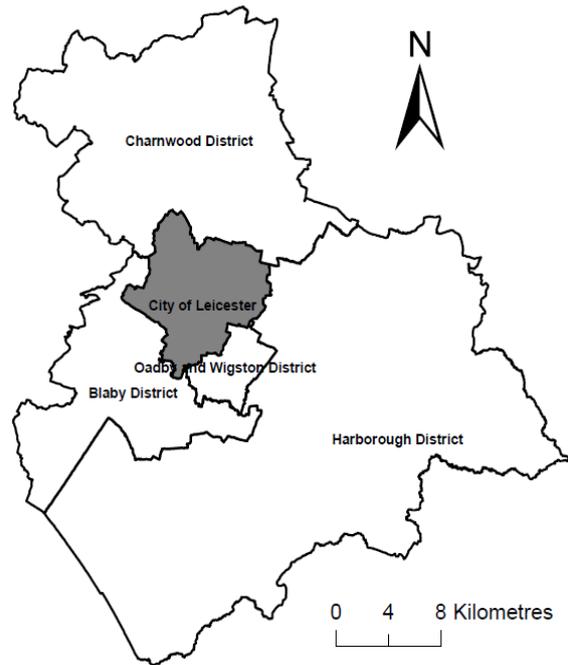
Fig 4: The Classified Leicester Image Derived from Landsat7 (ETM). The Digital Boundaries are © Crown Copyright and/or Database Right 2013. An Ordnance Survey/EDINA Supplied Service.



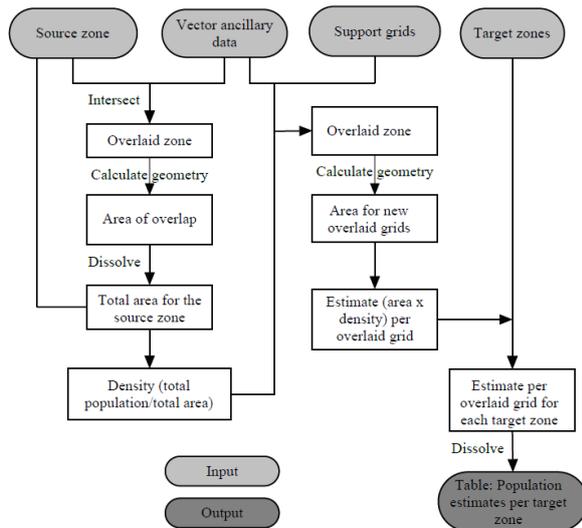
**Fig 5: A Binary Mask Derived from Land Cover Data Derived from Classified Landsat7 (ETM). The Digital Boundaries are © Crown Copyright and/or Database Right 2013. An Ordnance Survey/EDINA Supplied Service.**

**Table III Population totals for source zones used to implement pycnophylactic interpolation**

Unitary Authority	Total Population
Blaby District	90252
Charnwood District	153462
City of Leicester	279921
Harborough District	76559
Oadby and Wigston District	55795



**Fig 7: Source Zones Used for the Pycnophylactic Interpolation Method With the City of Leicester (Study Area) Shaded in Grey.**



**Fig 6: Implementation of the Binary Dasymetric Method**

### 3) Pycnophylactic Interpolation

Pycnophylactic interpolation can only be applied to two or more areas. In this analysis the aim was to derive an interpolated surface of the population count for Leicester over 100m and 30m grids which were then aggregated up to the MSOA, LSOA and OA target units. The unitary authority of Leicester was represented by a single polygon and so data for adjacent county districts (Charnwood, Harborough, Blaby, Oadby and Wigston) were included in the analysis to generate a pycnophylactic surface with five source zones (as in Figure 7). The total population for each source zone is shown in Table III.

The basic principle of the pycnophylactic interpolation is to create a smooth surface across the study area with no sudden changes across target zone boundaries, such that the total value of target region polygons must equal that of the source polygons with each source zone population being the same [12]. Figure 8 illustrates the general concept of the pycnophylactic interpolation. The method iteratively distributes populations, whilst seeking to smooth adjacent cells values and maintain total population volumes. It computes a continuous population density (per cell) in each source zone. The population density per cell is then smoothed repeatedly by replacing the value of each cell with the weighted average of its neighbours. The volume of the attributes within each source zone remains unchanged but varies smoothly at the boundaries.

The procedure for generating the pycnophylactic surface has been described by Qiu et al. [28] and involves the following steps: (1) converting the source zone data to raster grids; (2)

preserving the vector attributes in the raster; (3) computing the population density per grid cell; (4) calculating a new density by replacing the value of each cell with the weighted average of its neighbours; (5) estimating the density for each source zone using the new per cell density; (6) adjusting the new density by multiplying each cell value with the ratio between the original population and the estimated total population density of each source zone; (7) repeating steps 3-6 until no more adjustment is required for example when the maximum change in any pixel density values between iterations falls below a threshold level, such that zone total equals original value (the pycnophylactic condition); (8) obtaining the estimated interpolated gridded population of each target zone by summing the adjusted population density of each cell falling within each target zone. The implementation steps described above are illustrated in Figure 9.

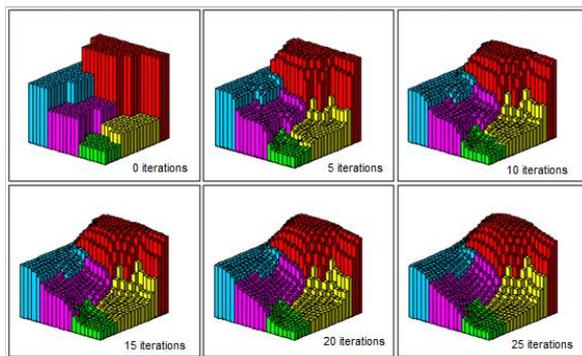


Fig 8: The Pycnophylactic Interpolation Method (Source: Tobler, W. R., 1992)

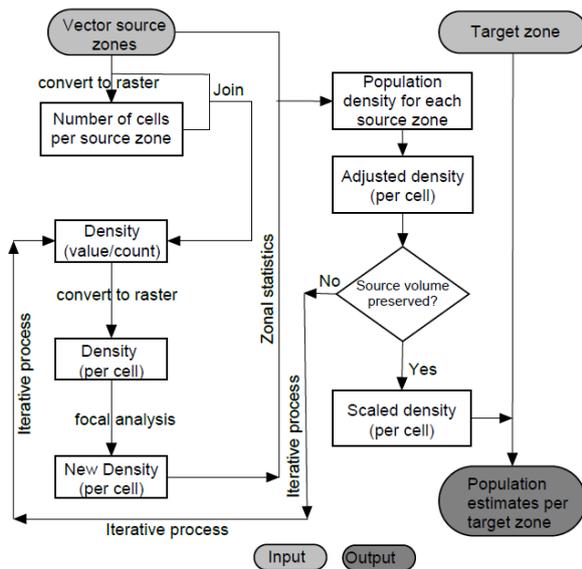


Fig 9: The Implementation Steps for the Pycnophylactic Interpolation Method

**D. Evaluation of surfaces**

The interpolated gridded pycnophylactic surfaces, the areal weighting and the dasymetric

population surfaces at 100m and 30m resolutions of the output grid were overlaid with the boundaries of MSOA, LSOA and OA target zones for Leicester and then aggregated to obtain estimates of the populations for these zones. These were assessed for accuracy by comparing the estimated populations with known census counts in each case. The boundaries of the target zones were spatially overlaid with the interpolated gridded population surfaces from which the target zone populations were calculated. Figure 10 shows an example of the results, in this case from the pycnophylactic interpolation at 100m resolution. The estimated target zone populations were then compared with known census counts in that target zone. The same procedure was repeated to obtain population estimates for the three U.K. census units, MSOAs, LSOAs and OAs that were used as the target zones.

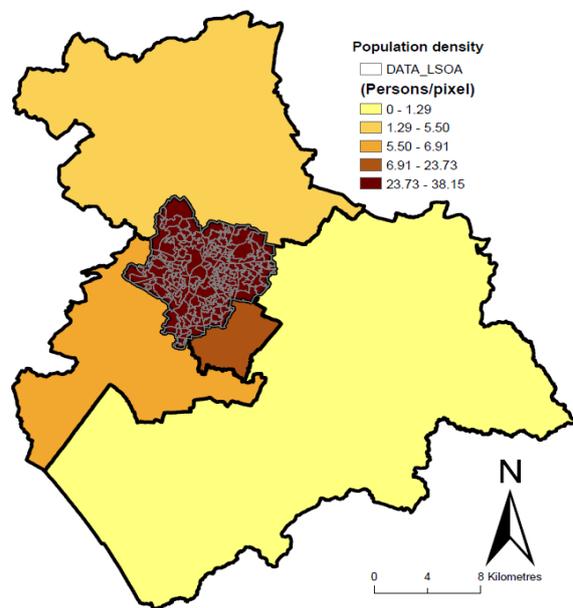


Fig 10: Leicester LSOAs Intersect Interpolated Gridded Pycnophylactic Surfaces at Resolutions of 100m Support Grids.

For each analysis residuals were calculated and mapped to visually explore the nature of the error [1,14,28,29]. The residual is calculated from the estimated population subtracted from the actual populations of each census unit. The accuracy of the interpolation is measured using the root mean squared error (RMSE) metric [1,14,23,29]. The RMSE metric gives a summary of the error within census units and was used to evaluate the different approaches. The error within a given source zone (RMSE) uses the absolute difference between estimated populations and the populations reported for the census units within each of the target zones and is calculated as in Equation 3.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - Y_i)^2} \quad [3]$$

Where:  $X_i$  is the known census count at zone  $i$ ,  $Y_i$  is the estimated population from the interpolation at zone  $i$ , and  $n$  is the number of target zones.

The RMSE metric has been found to be ‘less useful for comparing between different sets of source and target units’ [1, p.337], particularly where resolution change is involved. This is because the RMSE metric is affected by the absolute size of estimated values (e.g. MSOA counts are as expected larger than LSOA counts and would have a larger RMSE values). Previous research [1,14,29] has considered the variation in actual population of the target zones (e.g. MSOA and OA). To account for these variations, the RMSE score is divided by the average known population of each target zone to obtain the coefficient of variance (CoV). The CoV provides a relative error metric suitable for comparing values across the target zones. This is a useful metric as this research seeks to test performance over census areas of differing resolutions and CoV is more appropriate for cross-resolution comparisons. The CoV is calculated as in Equation 4.

$$CoV = \frac{RMSE}{\bar{x}} \quad [4]$$

Where:  $\bar{x}$  is the mean population of the known census count for each target zone.

#### IV. RESULTS AND DISCUSSIONS

The accuracy of the various interpolations was measured using the RMSE metric and CoV. Tables IV to IX summarise the performance measures for the various analyses in order of increasing accuracy. Tables IV and V show the interpolation results using 30m and 100m support grids respectively, aggregated at MSOA. Tables VI and VII show the interpolation results using 30m and 100m support grids respectively, aggregated at LSOA, and Tables VIII and IX show the interpolation results using 30m and 100m support grids respectively, aggregated at OA. The areal weighting method provides a baseline against which to compare the other techniques [1]. As expected, areal weighting performs least well in all the experiments undertaken.

Table IV shows the interpolation results using the 30m support grid, aggregated at MSOAs. The areal weighting method performed least well with RMSE score of 4486.9 and a CoV of 0.577. The pycnophylactic method slightly improves on areal weighting with a RMSE of 4233.4 and a CoV of 0.544. Interpolations using binary dasymetric with classified land cover data used as the ancillary data input are better than the pycnophylactic method. The

binary dasymetric methods using ancillary data input of differing spatial resolutions recorded slightly different CoV scores. The binary dasymetric model using land cover data derived from classified Landsat7 (ETM) 30m spatial resolution as the ancillary data input provided the best estimates among the methods tested with the lowest recorded RMSE of 2943.8 and a CoV of 0.379. The most striking feature in Table IV is that the binary dasymetric model using land cover data derived from classified resampled aerial photo data of 10m spatial resolutions as the ancillary data input achieved a RMSE values of 3314.4 which marginally improves to 3304.9 compared to a binary dasymetric model using land cover data derived from classified resampled aerial photo data of 3m spatial resolutions as the ancillary data input. This is surprising because higher resolution land cover data that offer greater spatial precision in the depiction of building locations does not automatically improves interpolation performance. One possible reason for this could be because they are both from the same source.

**Table IV Interpolation Results using the 30m Support Grids, Aggregated at MSOA**

Interpolation method	RMSE	CoV
Areal weighting using zone boundaries only	4486.9	0.577
Pycnophylactic interpolation	4233.4	0.544
Binary dasymetric using 10m resolution classified land cover	3314.4	0.426
Binary dasymetric using 3m resolution classified land cover	3304.9	0.425
Binary dasymetric using 30m resolution classified land cover	2943.8	0.379

Note: Mean population of target units is 7776.

Table V shows the interpolation results of the 100m support grids, aggregated to MSOAs in order of increasing accuracy. The results are similar to those in Table IV with the areal weighting method performing least well and the binary dasymetric model using land cover data derived from classified Landsat7 (ETM) 30m spatial resolution providing better target zone estimates. At MSOA, interpolations to 30m support grids are better compared to 100m support grids. Also, in contrast to Table IV, the RMSE value and CoV recorded for the binary dasymetric model using land cover data derived from 10m data marginally improves those recorded for the binary dasymetric model with 3m spatial ancillary data.

**Table V Interpolation Results Using the 100m Support Grids, Aggregated at MSOA**

Interpolation method	RMSE	CoV
Areal weighting using zone boundaries only	4934.1	0.635
Pycnophylactic interpolation	3974.9	0.511

Binary dasymetric using 3m resolution classified land cover	3668.7	0.472
Binary dasymetric using 10m resolution classified land cover	3661.5	0.471
Binary dasymetric using 30m resolution classified land cover	3579.7	0.460

Note: Mean population of target units is 7776.

Tables VI and VII show the interpolation results using 30m and 100m support grids, aggregated to LSOAs in order of increasing accuracy. The results recorded have similar pattern to those in Table V. The areal weighting method performed least well with RMSE of 1497.8 and 1805.3 (using 30m and 100m support grids respectively). The binary dasymetric model using 30m land cover data provided the best target zone estimates with a RMSE of 1087.5 and 1309.4 for the 30m and 100m support grids respectively.

**Table VI Interpolation Results Using the 30m Support Grids, Aggregated at LSOA**

Interpolation method	RMSE	CoV
Areal weighting using zone boundaries only	1497.8	1.001
Pycnophylactic interpolation	1368.5	0.914
Binary dasymetric using 3m resolution classified land cover	1173.9	0.784
Binary dasymetric using 10m resolution classified land cover	1155.6	0.772
Binary dasymetric using 30m resolution classified land cover	1087.5	0.726

Note: Mean population of target units is 1497.

**Table VII Interpolation Results using the 100m Support Grids, Aggregated at LSOA**

Interpolation method	RMSE	CoV
Areal weighting using zone boundaries only	1805.3	1.206
Pycnophylactic interpolation	1517.7	1.014
Binary dasymetric using 3m resolution classified land cover	1467.5	0.980
Binary dasymetric using 10m resolution classified land cover	1436.1	0.959
Binary dasymetric using 30m resolution classified land cover	1309.4	0.875

Note: Mean population of target units is 1497.

**Table VIII Interpolation Results Using the 30m Support Grids, Aggregated at OA**

Interpolation method	RMSE	CoV
Areal weighting using zone boundaries only	586.2	1.861
Pycnophylactic interpolation	516.8	1.641
Binary dasymetric using 3m resolution classified land cover	458.1	1.454
Binary dasymetric using 10m resolution classified land cover	447.4	1.420
Binary dasymetric using 30m resolution classified land cover	429.6	1.364

Note: Mean population of target units is 315.

Tables VIII and IX show the interpolation results for the 30m and 100m support grids respectively, aggregated to OAs in order of increasing accuracy. The results show a similar pattern to the LSOA results. The areal weighting method performed least well with RMSE of 586.2 and 761.9 (using 30m and 100m support grids respectively) and CoV of 1.861 and 2.419 (using 30m and 100m support grids respectively). The binary dasymetric model using 30m land cover provided the best target zone estimates at this resolution of interpolation with a RMSE of 429.6 and 503.5 (for 30m and 100m support grids respectively) and CoV of 1.364 and 1.598 (for 30m and 100m support grids respectively).

**Table IX Interpolation Results using the 100m Support Grids, Aggregated at OA**

Interpolation method	RMSE	CoV
Areal weighting using zone boundaries only	761.9	2.419
Pycnophylactic interpolation	664.3	2.109
Binary dasymetric using 3m resolution classified land cover	630.1	2.000
Binary dasymetric using 10m resolution classified land cover	614.4	1.950
Binary dasymetric using 30m resolution classified land cover	503.5	1.598

Note: Mean population of target units is 315.

The binary dasymetric method shows significant improvement when compared with the pycnophylactic and areal weighting methods: it provides the best estimates among the models tested. This improved performance, quantified in this research, is because the technique uses land cover data to constrain the population distribution to only populated areas. The results presented show how the underlying assumptions of each interpolation technique and the scales of analysis interact to influence the target zone estimates. The dasymetric method was found to consistently provide better target zone estimates when compared to other areal interpolation techniques.

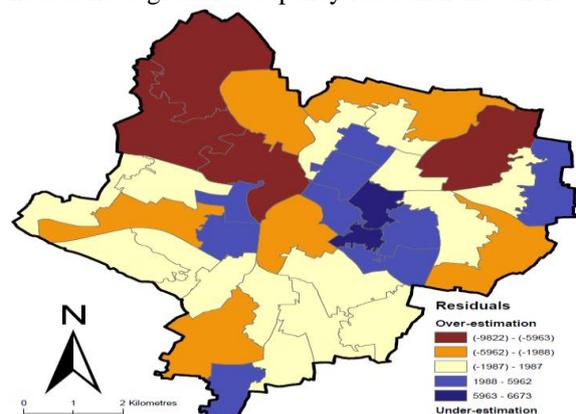
The binary dasymetric model using land cover data derived from classified Landsat7 (ETM) 30m spatial resolution as the ancillary data input provided the lowest recorded RMSE score for all the models tested, for the three target zones compared to land cover data derived from classified resampled aerial photo of 10m and 3m spatial resolutions. The expectation is that high resolution satellite image, which appears to offer greater spatial precision in identifying urban extent, could lead to reduction in land cover classification error. This is because the spectral signatures for each land cover type are likely to generate as little confusion as possible with a clear separation of land cover classes before classification. A possible reason for this result is the classification algorithm used in this study, the maximum likelihood

classifier. Maximum likelihood classification algorithm can provide reasonably good classification results for Landsat imagery [30,31]. The algorithm classifies land cover based on spectral signatures at per pixel level, while ignoring spatial features in an image. However, there are a number of issues related to using maximum likelihood classifier for medium and high resolution imagery. This is because a significant proportion of medium and high spatial resolution imagery in urban areas can be affected by shadows [32]. In this study, extracting urban land cover from resampled aerial photo data was more difficult compared to using the Landsat (ETM) source. Lu et al. [33] have shown how the use of spatial features improves land cover classification, especially when high spatial resolution images are used. Object-based classification provides an alternative for classifying remotely-sensed images into thematic map. Lu et al. [30] compared object-based classification with maximum likelihood and found object-based classification to be especially valuable for higher spatial resolution images. The object-based classification algorithm was not applied in this study. Also, the performance of 10m and 3m resampled aerial photo data can be attributed to using land cover information of different resolutions of the same source.

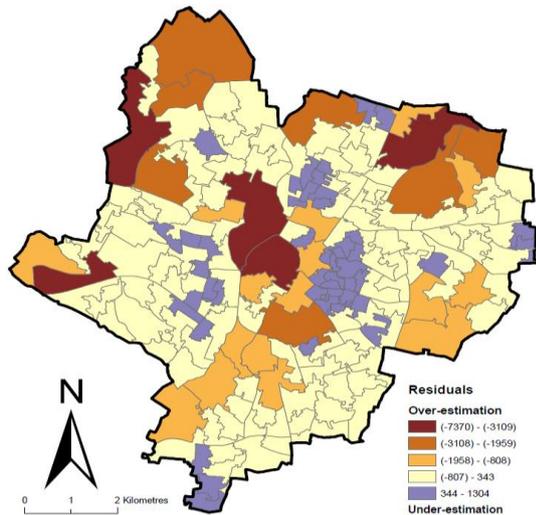
The interpolation results aggregated at MSOAs have larger RMSE values than those aggregated at LSOAs, which are also larger than those aggregated at OAs. This is expected because the RMSE metric is affected by the absolute size of estimated values and the target size population for an OA is less than that of an LSOA, which is also less than that of an MSOA. The CoV scores, which are appropriate for comparison across target zones, show interpolation results from Leicester unitary authority to 30m support grids, aggregated at MSOAs provided the lowest CoV score among the solutions tested and for the three census areas used as the target zones. This indicates that larger census units are more likely to produce better results as it shows improvements in RMSE and the values of CoV as the size of the spatial aggregation increases. This result is not surprising as one would expect higher accuracies when values are disaggregated over coarser spatial units. This result is similar to findings of Comber et al. [34] where a combination of pycnophylactic interpolation with the dasymetric method was used to create the National Agricultural Land Use Dataset. They reported improvement in  $R^2$  and RMSE values for Arable and Grass land uses for Kent, U.K. as the size of the spatial aggregation increases by the plots from 1 km<sup>2</sup> to 25 km<sup>2</sup>.

The residuals in all the census units tested were calculated and mapped to show the spatial distribution of the error (Figures 11, 12 and 13). The class intervals are shown by standard deviation from

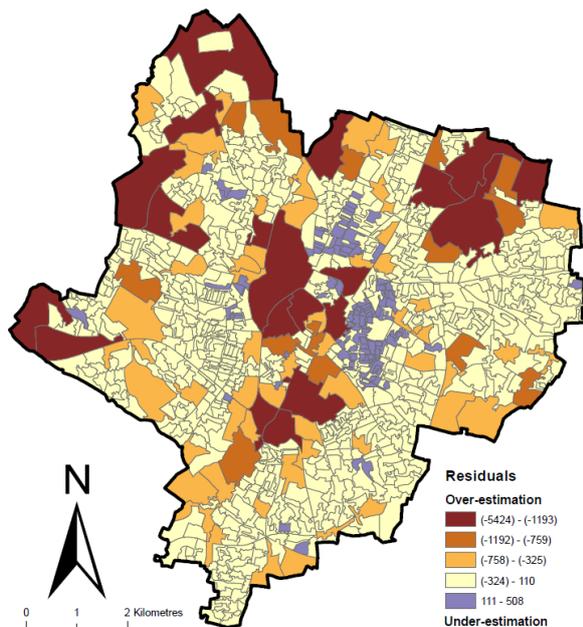
the mean error for each target zone. Standard deviations are the best way to symbolise normally distributed quantitative data on maps, as it makes classes easier to interpret. From the residual maps presented in Figures 11 to 13, some patterns persist across scales. It becomes increasingly clear that a degree of spatial ‘smoothing’ is present in the estimates. That is, the very densely populated inner city OAs are underestimated, the less dense band running alongside the river north-south through Leicester is overestimated, and many large rural OAs are overestimated. Evidently, the residual maps show more census areas are subject to overestimation, as compared to underestimation, at greater than one standard deviation. The residual maps show that relatively large rural census units tend to be overestimated while relatively small urban census units tend to be underestimated. This is because they are designed to have a common target population count [3]. Similar patterns have been found by other researchers e.g. [14,29], where relatively large rural blocks tend to be overestimated while relatively small urban blocks tend to be underestimated. The underestimated units are mainly the smaller census areas in the more densely populated areas such as the city centre while the overestimated units are the larger spatial areas in the less densely populated areas away from the city centre. A possible reason for this is that the satellite data being used as the ancillary data input may be more likely to identify houses and other built-up areas but not how many people live inside them. It is also likely in some areas there may be socioeconomic or cultural reasons why some houses have greater occupancy rates than the others.



**Fig 31: The spatial Distribution of Residuals at MSOA from a 100m Gridded Pycnophylactic Population Surface. The Mean Count Error is 0 and a Standard Deviation of 3975. The Digital Boundaries are © Crown Copyright and/or Database right 2013. An Ordnance Survey/EDINA Supplied Service.**



**Fig 42: The spatial distribution of residuals at LSOA from a 30m gridded dasymetric population surface using land cover data derived from classified resampled aerial photo data of 3m spatial resolutions as the ancillary data input. The mean count error is -233 and a standard deviation of 1151. The digital boundaries are © Crown Copyright and/or database right 2013. An Ordnance Survey/EDINA supplied service.**



**Fig 53: The spatial distribution of residuals at OA from a 30m Gridded Dasymetric Population Surface using Land Cover Data Derived from Classified Resampled Aerial Photo Data of 10m Spatial Resolutions as the Ancillary Data Input. The Mean Count Error is -107 and a Standard Deviation of 434. The Digital Boundaries are © Crown Copyright and/or Database right 2013. An Ordnance Survey/EDINA Supplied Service.**

## V. CONCLUSION

This study has developed a novel and comprehensive analysis of the operation of three classic spatial interpolation approaches and how they

interact with different target zone sizes, support grids and different scales of ancillary data. The results show how the underlying assumptions of each interpolation technique and the scales of analysis interact to influence the target zone estimates. The dasymetric method was found to consistently provide better target zone estimates when compared to other areal interpolation techniques. Much previous research using dasymetric methods have used land cover information derived from classified satellite imagery as the ancillary data input [1,7,15]. However, deriving such information from satellite imagery requires specialized skills and such data cannot determine population density levels, providing a potential source of error associated in analyses using such data as the ancillary data input. This study evaluated land cover data classified at different spatial resolutions (30m, 10m and 3m) as ancillary data input to the dasymetric method and found the coarsest resolution data generated the best results, with the lowest values of RMSE and CoV for all the models tested. These results, along with the free availability of 30m spatial resolution remote sensing data (Landsat etc.) and the ease with which it can be classified into urban and non-urban areas suggests its suitability as input for the dasymetric method. Thus this research suggests that additional costs and computational effort associated with finer scale remote sensing imagery (e.g. 10m and 3m resolution) has no analytical advantage and that the quality of the land cover data is not as important as its ability to predict the population estimates.

This study provides an important contribution to knowledge, with respect to estimating population surfaces. Fine scale estimates of spatial population have relevance for a broad range of applications, and therefore the findings of this research are of value beyond the field of Geographical Information Science. Research in the field of small area population estimates remains relevant because of the absence of a universally accepted methodology in estimating population surfaces. There is the need to apply areal interpolation techniques to different areas to be able to understand why a particular method consistently provides better target zone estimates.

## ACKNOWLEDGEMENT

This work was undertaken as part of a PhD funded by the Petroleum Technology Development Fund (PTDF) under the Federal Government of Nigeria [PTDF/E/OSS/PHD/MIJ/316/10]. The authors would like to express gratitude to Ordnance Survey, U.K. for providing 25cm Ortho-rectified aerial photograph covering Leicester area.

## REFERENCES

- [1] Langford, M. An Evaluation of Small Area Population Estimation Techniques Using Open Access Ancillary Data. *Geographical Information Science* 2013, 45, 324-344.
- [2] Harris, R. J.; Longley, P. A. Creating small area measures of urban deprivation. *Environment and Planning A* 2002, 34, 1073-1093.
- [3] Martin, D. From enumeration districts to output areas: Experiments in the automated creation of a census output geography. *Population Trends* 1997, 88, 36-42.
- [4] Martin, D. Optimizing census geography: The separation of collection and output geographies. *International Journal of Geographical Information Science* 1998, 12, 673-685.
- [5] Leyk, S.; Nagle, N. N.; Buttenfield, B. P. Maximum Entropy Dasymetric Modeling for Demographic Small Area Estimation. *Geographical Analysis* 2013, 45, 285-306.
- [6] Goodchild, M. F.; Lam, N. Areal interpolation: A variant of the traditional spatial problem. *Geo-Processing* 1980, 1, 297-312.
- [7] Zandbergen, P. A.; Ignizio, D. A. Comparison of Dasymetric Mapping Techniques for Small-Area Population Estimates. *Cartography and Geographic Information Science* 2010, 37, 199-214.
- [8] Wu, S.; Qiu, X.; Wang, L. Population estimation methods in GIS and remote sensing: a review. *Geographic Information Science and Remote Sensing* 2005, 42, 80-96.
- [9] Lam, N. S. Spatial Interpolation Methods: A Review. *The American Cartographer* 1983, 10, 129-149.
- [10] Langford, M.; Higgs, G. Measuring Potential Access to Primary Healthcare Services: The Influence of Alternative Spatial Representations of Population. *The Professional Geographer* 2006, 58, 294-306.
- [11] Hewko, J.; Smoyer-Tomic, K. E.; Hodgson, M. J. Measuring neighbourhood spatial accessibility to urban amenities: does aggregation error matter? *Environment and Planning A* 2002, 34, 1185-1206.
- [12] Tobler, W. Smooth Pycnophylactic Interpolation for Geographical Regions. *Journal of the American Statistical Association* 1979, 74, 519-530.
- [13] Langford, M.; Maguire, D. J.; Unwin, D. J. The area transform problem: Estimating population using remote sensing in a GIS framework. In: Masser, I. and Blakemore, M. (eds) *Handling Geographical Information: Methodology and Potential Applications*. London: Longman 1991, 55-77.
- [14] Eicher, C.; Brewer, C. Dasymetric mapping and areal interpolation: Implementation and evaluation. *Cartography and Geographic Information Science* 2001, 28, 125-138.
- [15] Mennis, J. Generating surface models of population using dasymetric mapping. *The Professional Geographer* 2003, 55, 31-42.
- [16] Green, N. E. Aerial photographic analysis of residential neighbourhoods: An evaluation of data accuracy. *Social Forces* 1956, 35, 142-147.
- [17] Lo, C. P. Population Estimation Using Geographically Weighted Regression. *GIScience and Remote Sensing* 2008, 45, 131-148.
- [18] Lillesand, T. M.; Kiefer, R. W., Eds.; *In Remote sensing and image interpretation*; John Wiley and Sons.: New York, 1987.
- [19] Lo, C. P. Estimating Population and Census Data. *American Society for Photogrammetry and Remote Sensing* 2006, 337-377.
- [20] Langford, M.; Unwin, D. Generating and mapping population density surfaces within a geographical information system. *The Cartographic Journal* 1994, 31, 21-26.
- [21] Cromley, R. G.; Hanink, D. M.; Bentley, G. C. A Quantile Regression Approach to Areal Interpolation. *Annals of the Association of American Geographers* 2011, 102, 763-777.
- [22] Mennis, J. Dasymetric Mapping for Estimating Population in Small Areas. *Geography Compass* 2009, 3, 727-745.
- [23] Tapp, A. F. Areal Interpolation and Dasymetric Mapping Methods Using Local Ancillary Data Sources. *Cartography and Geographic Information Science* 2010, 37, 215-228.
- [24] Su, M.; Lin, M.; Hsieh, H.; Tsai, B.; Lin, C. Multi-Layer Multi-class Dasymetric Mapping to Estimate Population Distribution. *Science of the Total Environment* 2010, 408, 4807-4816.
- [25] Langford, M. Obtaining population estimates in non-census reporting zones: An evaluation of the 3-class dasymetric method. *Computers, Environment and Urban Systems* 2006, 30, 161-180.
- [26] Qiu, F.; Cromley, R. Areal Interpolation and Dasymetric Modeling. *Geographical Analysis* 2013, 45, 213-215.
- [27] Sridharan, H.; Qiu, F. A Spatially Disaggregated Areal Interpolation Model Using Light Detection and Ranging-Derived Building Volumes. *Geographical Analysis* 2013, 45, 238-258.
- [28] Qiu, F.; Zhang, C.; Zhou, Y. The Development of an Areal Interpolation ArcGIS Extension and a Comparative Study. *GIScience and Remote Sensing* 2012, 49, 644-663.
- [29] Mennis, J.; Hultgren, T. Intelligent dasymetric mapping and its application to areal interpolation. *Cartography and Geographic Information Science* 2006, 33, 179-194.
- [30] Lu, D.; Li, G.; Moran, E.; Freitas, C. C.; Dutra, L.; Anna, S. J. S. In *In A Comparison of Maximum Likelihood Classifier and Object-Based Method Based on Multiple Sensor Datasets for Land-Use/Cover Classification in the Brazilian Amazon*; Proceedings of the 4<sup>th</sup> GEOBIA; Rio de Janeiro, Brazil, 2012; pp 20.
- [31] Blaschke, T. Object based image analysis for remote sensing. *ISPRS Journal of Photogrammetry and Remote Sensing* 2010, 65, 2-16.
- [32] Zhou, W.; Huang, G.; Troy, A.; Cadenasso, M. L. Object-based land cover classification of shaded areas in high spatial resolution imagery of urban areas: A comparison study. *Remote Sensing of Environment* 2009, 113, 1769-1777.
- [33] Lu, D.; Hetrick, S.; Moran, E. Land cover classification in a complex urban-rural landscape with QuickBird imagery. *Photogrammetric Engineering and Remote Sensing* 2010, 76, 1159-1168.
- [34] Comber, A.; Proctor, C.; Anthony, S. The Creation of a National Agricultural Land Use Dataset: Combining Pycnophylactic Interpolation with Dasymetric mapping technique. *Transactions in GIS* 2008, 12, 775-791.