*Article*

# Population Estimates from Orbital Data of Medium Spatial Resolution: Applications for a Brazilian Municipality

**Járvis Campos [1]**, **José Irineu Rangel Rigotti [2]**, **Emerson Augusto Baptista [3],***,
**Antônio Miguel Vieira Monteiro [4]** and **Ilka Afonso Reis [5]**

1   Departamento de Demografia e Ciências Atuariais (DDCA),
    Universidade Federal do Rio Grande do Norte (UFRN), Natal 59078-970, Brazil; jarvis@ccet.ufrn.br
2   Department of Demography, Universidade Federal de Minas Gerais (UFMG),
    Belo Horizonte 31270-901, Brazil; rigotti@cedeplar.ufmg.br
3   Asian Demographic Research Institute (ADRI), Shanghai University, Shanghai 200444, China
4   LiSS/CGOBT, Instituto Nacional de Pesquisas Espaciais, São José dos Campos 12227-010, Brazil;
    miguel.monteiro@inpe.br
5   ICEX, Universidade Federal de Minas Gerais (UFMG), Belo Horizonte 31270-901, Brazil; ilka@est.ufmg.br
*   Correspondence: emersonaug@gmail.com

check for updates

**Abstract:** In recent decades, there has been an increase in the search for more detailed information on population dynamics, given the growing demand for more sustainable economic, social, and environmental planning. The dissemination of Geographic Information Systems (GIS) has contributed to the development of methodologies for the field of population estimates for small areas. To support more sustainable policies, this study aims to evaluate the capacity and contribution of the orbital images (Landsat ETM+) for the production of post-census population estimates for the municipality of Contagem, Minas Gerais, Brazil. Firstly, models were built using the average of the reflectance of the spectral bands of the Landsat 7 ETM+ for each special intra-municipal unit, called the census sector, as explanatory variables for the population density. Secondly, this study constructed models that use the reflectance and the distributed population at the level of the pixels of the images. All models were tested through internal validation procedures, external validation, and comparative analyses with post-census estimates. Internal validation presented excellent results (below 7%), while in external validation, the method at the level of the pixels presented consistent results, below 1% relative error. These results provide useful clues and can help policymakers in the development of more sustainable and effective public policies, insofar as population estimates are extremely important for the planning of any society.

**Keywords:** remote sensing; dasymetric mapping; population estimates; small areas; Brazil; sustainable development

## 1. Introduction

Information on the size, composition, and pace of population change has been the subject of increasing interest in the public sector, for economic policy and for policy-making in other areas, including health, education, sustainability, and the environment. To date, the most complete and reliable type of source for data on the population of countries and their geographical subdivisions has been a census based on household interviews.

If, on the one hand, the periodic enumeration of censuses in many countries allows the capture of population dynamics, the time interval (every 5 or 10 years) is usually not adequate to the pace of

population change. Moreover, census information is limited due to high costs, despite the growing demand for more detailed information, as the social and economic demands on governments become increasingly complex [1–5].

Alternative sources of information are used, such as sample surveys (although equally expensive), civil records (especially for the enumeration of vital statistics), school censuses, and various government registries, although this information is not available for most countries (or its administrative subdivisions). Added to this, there are problems with underreporting and the quality of information is sometimes lacking [1–3].

Therefore, a set of methods that use symptomatic variables has been developed for the production of estimates of small areas; these approximations of population change are fundamental for planning in various sectors [2–7]. A symptomatic variable is an auxiliary variable strongly correlated with the population growth of a given locality, and one that changes in time in accordance with changes in the population volume. It is usually used for the production of estimates in smaller areas [4,8], but also for estimates of some large areas.

For some time, the use of data derived from space programs, in particular satellite imagery data, has drawn attention to the possibility of its use in population estimates [8,9]. The dissemination and intensification in the use of Geographic Information Systems (GIS), the development of spatial analysis techniques, and the growth of the availability of remote sensing images have contributed to the development of a set of methods with very promising applications for population estimates of small areas [10–19].

There are several elements in the area of remote sensing, in addition to spatial resolution, such as spectral and/or radiometric images, which can interfere with the construction of reliable estimates. Some examples are the level of cloud cover, the procedures adopted in the image processing and classification steps, in addition to the assumptions and limitations of the dasymetric method (which is a population distribution method) and the model used for the calculation of the estimates. However, spatial resolution is a key component in the production of the estimates, and the choice of the image (with a given spatial resolution) determines, to a large extent, the processing and classification methods, the statistical model to be adopted and, in the end, the quality and accuracy of the estimates produced [10,20].

In recent years, advances in remote sensors and in the dissemination of high spatial resolution images have expanded the studies and applications of remote sensing for the calculation of estimates [21–23]. However, most high spatial resolution images (with pixel resolution from approximately 5 meters) are not available for free. Among those that are freely available, there is no historical series that includes the period of two censuses, a fundamental element for assessment of the quality of the estimates produced. Roughly speaking, the parameters of the models are defined in the first census, while in the next census, the models are used to calculate the estimates, based on comparison with the data generated from the last census.

Another limitation of high-resolution images is their limited scope of application to more extensive areas, since a smaller area for the field of view for each scene makes the construction of large mosaics practically unviable in relation to the quantity and processing time of the images (except for large private sector companies, such as Google, Amazon, etc.). In the public sector, the systematic production of estimates for large areas (states / country), as well as the cost and processing time of high-resolution images, limit the production of official estimates. It is true that this statement is dated, as technological advances tend, at an accelerated pace, to overcome processing restrictions. However, the demand for the production of estimates of small areas, replicable to large areas, justifies the choice of free-access images, with a historical series available and consolidated in relation to the quality of the calculated estimates.

In the direction of these advances, this paper aims to contribute to the evaluation of the capacity of orbital images of medium spatial resolution (Landsat ETM+) in the production of post-census population estimates at the municipal level, as a way to provide useful clues that can help policymakers in the development of more sustainable and effective policies. The choice of the Landsat ETM+ satellite

sensor is justified because its images have relatively good spatial resolution, among the satellites that offer free images and have time series available.

Remote sensing estimates were calculated for the municipality of Contagem (Minas Gerais, Brazil) in 2000, 2010, and 2015. Over the decades, there has been an evolution in the methods used to produce estimates via Remote Sensing. The first methods considered the direct relationship between growth in urban areas and population size [9,24–26]. In order to advance studies that relate population to the expansion of urban areas, methods have been developed that analyze the correlation between population and different types of land use [12,27–30]. There is also a set of methods that estimate the population through the product of the number of housing units and the number of people who normally live in these units [21–23]. In these methods, the population density was associated with characteristics or land-use classes. More recently, a set of pixel-based methods has emerged, in which the population can be directly correlated to the spectral reflectance of the pixel image values [10].

In general, the data structure consists of the construction of regression models, based on census data and orbital images, for the years 2000 and 2010. In these models, the dependent variable is the population, while the explanatory variables are the reflectance of the bands of the Landsat images, in addition to some additional variables that have been tested in relation to their ability to predict, such as whether the area is urban, rural or slums. The models constructed from 2000 and 2010 data were calibrated using 2/3 of the pixels of the Landsat images. The other 1/3 of the pixels were selected randomly for the test sample, that is, for the analysis of the error of the models in relation to the population size prediction capacity, a step that was called "internal validation".

Secondly, the estimates calculated for 2010, based on the models constructed with data from 2000, were compared with the population listed in the 2010 census, a step called "external validation". Finally, models constructed from 2000 and 2010 data were used to prepare estimates for the year 2015, which, in turn, were compared with other post-census estimates, including a projection by the Brazilian Institute of Geography and Statistics (IBGE) and estimates produced via three other demographic methods that use symptomatic variables (simple extrapolation, vital rates and census ratios). In this way, we hope to contribute to the evaluation of the capacity of orbital images of medium spatial resolution for the production of post-census estimates at the municipal level.

## 2. Materials and Methods

Contagem is a municipality in the state of Minas Gerais, located in the southeast region of the country. It is the third most populous municipality in the state (31st most populous in the country, out of a total of 5,570), with an estimated population of 659,070 inhabitants in 2018, a population density of 3090.33 hab/km$^2$ in 2010, and an annual growth rate of 1.15%. It is a metropolitan municipality with very peculiar occupation characteristics that justify its selection. Contagem is characterized by the presence of extensive vertical areas with high population density, as well as neighborhoods with less verticalization and lower density, in addition to an important industrial complex.

One of the biggest problems in producing population estimates in urban areas is industrial complexes. Their occupation profile is marked by scattered residences, but they are normally classified as having a high degree of urbanization (that is, high population density) due to the dense road networks and extensive impermeable areas, which result in spectral behaviors consistent with areas of intense occupation [29]. The heterogeneous characteristics of Contagem, therefore, justify its choice, due to the challenge of the classification stage (largely due to the presence of industry), and due to the type of growth observed in the municipality, marked both by the expansion of the urban area, and by verticalization. In addition, the selection of a high-density site represents the most explored municipality profile in the literature on remote sensing estimates.

In general, the construction of estimates through remote sensing begins with census years. Through the use of auxiliary data (satellite imagery), the geographic distribution of population data with a high level of spatial detail obtained in the census (census tracts) is transformed into a spatial distribution of the population that is even more refined. This procedure is called a surface or dasymetric

mapping. After performing the dasymetric mapping for a given census year, an econometric model is constructed, based on the relationship between population (dependent variable) and reflectance (explanatory variables), which is information contained in the pixels of the satellite images of the occupied areas, as defined by dasymetric mapping. The econometric model is then applied to another dasymetric map corresponding to a post-census year, and the estimate for that year is thereby obtained.

### 2.1. Dasymetric Mapping

In this study, two geographic units were used to construct the estimates: census tracts and statistical grids cells. Census tracts correspond to the smallest administrative unit of the Brazilian demographic census (of IBGE) in which information on population size is available. In 2010, the municipality of Contagem had 884 census tracts. The statistical grids cells, in turn, is a dasymetric mapping application provided by IBGE, which starts with the use of satellite images to represent population size data in grids cells, independently of the official administrative divisions. Thus, statistical grids cells make it possible to analyze data in small geographical units. They also solve the problem of area incompatibility (MAUP) over time [31]. The grids cells have a resolution of 200 meters for urban areas and 1 km for rural areas. In 2010, Contagem presented 4462 cells (statistical grids cells).

Dasymetric mapping, in turn, consists of taking the source unit population totals (census tracts and statistical grids cells, in the case of this paper) and distributing them to target unit grid cells (i.e., pixels) of some defined spatial resolution, where the disaggregation may be informed by some data detected remotely. The dasymetric process consists, therefore, in census counts tied to some GIS-administrative boundary at a given level of detail, depending on the location and year. Those counts can then be disaggregated based on remote sensing information [10,15,19,28–30,32–39]. These models allow the discrimination between occupied and unoccupied physical space structures (such as vegetation and rivers, among others), although it is important to observe the uncertainties inherent in this process [29,33,36,38]. Important programs monitor the special distribution of the population in grids cells, such as the Gridded Population of the World (GPW) project at Columbia University [18], the Urban Atlas project [34,35] by the European Environment Agency, and the WorldPop project at the University of Southampton [40].

For the construction of models and estimates, the first step was the selection of Landsat ETM+ satellite images for the years 2000, 2010 and 2015. Landsat imagery [41] is freely available, and with a spatial resolution of 30 meters, it has a large field of view, which reduces processing time and enables it to be used for more extensive areas. The imagery includes a long historical series (for example, including the period of the last two Brazilian censuses), which enables the construction and validation of models for the calculation of estimates. Despite the existence of a large international literature on the use of Landsat imagery for the calculation of estimates for small areas [8,11,13,23,25,27,30,32,33,42,43], in Brazil there are still few studies that use this imagery for this type of application (exceptions are the works of [21,26,44–46], which justify the choice of this type of image with medium spatial resolution).

The Landsat ETM+ satellite images are composed of eight spectral bands, six bands with a spatial resolution of 30 meters, a panchromatic band with a resolution of 15 meters, and the thermal band with a resolution of 60 meters. Table 1 shows the main characteristics of the ETM + sensor, with emphasis on the spectrum bands (defined for each band).

The digital numbers of the satellite images correspond to the radiance, for each band of the image. However, the regression models used to calculate the estimates consider the surface reflectance. The United States Geological Survey (USGS) has a project, called ESPA, which provides satellite images with atmospheric correction and with the surface reflectance data in the pixels. The ESPA project eliminated the need to calculate reflectance for the construction of regression models. In addition, the images of the ESPA project have geometric correction; that is, they are registered and made available in "GeoTiff" format, in the WGS84 system [47].

**Table 1.** Main characteristics of the ETM+ Sensor.

| Characteristics | Sensor Parameters |
|---|---|
| Spectral bands (μm) | Band 1 - 0.45 a 0.52 |
| | Band 2 - 0.53 a 0.61 |
| | Band 3 - 0.63 a 0.69 |
| | Band 4 - 0.78 a 0.90 |
| | Band 5 - 1.55 a 1.75 |
| | Band 6 - 10.4 a 12.50 |
| | Band 7 - 2.09 a 2.35 |
| | Band 8 - 0.52 a 0.90 |
| Spatial resolution | 15 meters (panchromatic band) |
| | 30 meters (bands 1 to 5 and 7) |
| | 60 meters (band 6) |
| Radiometric resolution | 8 bits (256 gray levels) |
| Size of the scenes | 170 km (north-south)/183 km (east-west) |

Source: USGS, 2017.

In order to produce estimates for the municipality of Contagem (Minas Gerais, Brazil), images were selected for 2000, 2010, and 2015 with a percentage of clouds below 10%, and with dates that are closest to the reference date of the 2000 and 2010 censuses. The Landsat 7 scenes collected from 30 May, 2003 have data gaps due to the failure of the Scan Line Corrector (SLC). Therefore, there is loss of information in 22% of the pixels of all images from that date. In order to use a single sensor for the period under analysis, this study used a filter procedure (pixel interpolation) for the correction of pixels with no information.

The next step consisted of the classification of Landsat ETM+ satellite imagery (recorded, with no atmospheric correction and with low cloud percentage), in which the occupied area of Contagem was mapped for the three years under analysis (2000, 2010, and 2015). For this, the Maxver supervised classification method was applied to 4 classes (occupied areas, vegetation, water and soil) and the 3 selected images, from the SCP plugin of the QGIS software. In the Maxver method, a set of training samples is selected for each of the classes previously defined. In order to determine which class a pixel belongs to, the Maximum Likelihood (Maxver) classification assesses the probability that a given class $w_i$ is the correct one for a pixel x (where: M is the total number of classes):

$$p(w_i|\overline{x}), i = 1, \ldots, M. \tag{1}$$

$$\overline{x} \in w_{i \text{ if}} \Rightarrow p\left(w_i|\overline{x} > p(w_j|\overline{x})\right)_{\text{for all}} i \neq j \tag{2}$$

From the training samples, it is possible to estimate the probability distribution of each class, and each pixel is assigned to the class with the highest probability [20]. For example, a limit of 95% means that 5% of pixels, the least likely, will be ignored. This procedure aims to compensate for the possibility of errors in the training phase, which tend to increase the overlap between the probability distributions of the classes, or even pixels that are at the limit between two classes. In addition to the certain degree of arbitrariness in the selection of training samples, the automatic classification method used in medium resolution images presents classification errors between exposed soil cover and areas of human occupation, as well as relief problems (shadows generated), which in many cases makes it difficult to define urban areas. Further, the variety of coverage existing in urban areas (such as concrete, asphalt, roofs, vegetation within urban areas, among others) makes the distinction by automatic classification complex, due to the various spectral responses existing in the occupied areas. In view of the importance of qualitative analysis for a good classification, this study adopted the hybrid method, which consists of automatic classification followed by the visual interpretation of images with the aim of improving the distinction between occupied and unoccupied areas, and, consequently, the estimates produced. With this procedure to minimize problems related to the classification stage, it was possible to analyze the results of the estimates in light of the advantages and limitations of

the proposed model, and, mainly, in light of the different occupation characteristics existing in the municipality of Contagem. Then, a second classification step was carried out, which consisted in the manual interpretation of high spatial resolution satellite imagery (Google) for the areas and years studied, representing an important quality gain in the delimitation of areas of human occupation.

After classification of the occupied areas, the next step consisted of the dasymetric mapping of the population, which was based on the distribution of the population on the same spatial scale as the Landsat ETM+ satellite images (pixels with a resolution of 30 meters). With the distribution of the Contagem population on the pixel scale—by the criterion of homogeneous distribution for each set of pixels located in a given census tract and in the statistical grids cells—georeferenced databases were constructed and made compatible between the population information and the reflectance of each spectral band of Landsat imagery (bands 1 to 5 and 7) by means of the use of vector grids cells. This information served as a basis for building the models in the next step.

## 2.2. Methods by Zones and Pixels

### 2.2.1. Method by Zones

Harvey's works [48,49], the most cited in international literature, are considered classics among those that deal with the correlation between population and the pixels of satellite imagery. Therefore, this study used the methodologies proposed by him and replicated by [44,46], which deal with this relationship. Nevertheless, several other methods (with different satellite imagery of different resolutions) were used for the production of estimates through remote sensing [8,11,13,16,17,21–23,25,26,28,29,33,43,48,49].

Harvey [48] used images from the TM sensor of the Landsat satellite to estimate the population of two districts in Australia: Ballarat statistical district, which was used for the construction of regression models, and the Geelong statistical district, used for the external validation of the models. Starting from the pixels belonging to the same intra-urban spatial unit (individual zones), the average values of the reflectance of bands 1 to 5 and 7 of the Landsat imagery (for each individual zone) were used as explanatory variables in a linear regression model. Ordinary Least Squares Models, with normal error distribution, were used in this study:

$$p_i = \beta_0 + \sum_{j=1}^{k} \beta_j r_{ij} + \varepsilon_i \qquad (3)$$

In the model above, $p_i$ represents the population density of census tract $i$, $r_{ij}$ is the average reflectance of the pixels of census tract $i$ in the $j$-th sensor band, β0, and βj, j = 1, 2, ... k, are the parameters to be estimated and $\varepsilon_i$ represents the part of the population density of the census tracts that is not explained by the regression model.

Each explanatory variable usually represents a band of the satellite imagery, being the average of the reflectance of the pixels inserted in a given spatial unit (called individual zones by the author). On the other hand, the dependent variable can be transformed to improve model estimates (such as population transformation to logarithmic function) [50].

After internal validation of the regression models, from the selection of training samples for the same district with which the model was created, the regressions were analyzed from 132 spatial units. The dependent variable chosen was the mean of the population density, while the explanatory variables were the means of the reflectance of the pixels belonging to the same spatial unit [48]. The results showed that the application of the models to the Geelong district reached $R^2 = 0.84$ and the linear correlation coefficient was equal to 0.92 between the estimated and actual values of the population density, while the median error was 17.4% in the training sample and 18.4% in the external validation. On the other hand, the urban population obtained errors of 1% and -3%, to the training sample and external validation, respectively (application to the Geelong district).

### 2.2.2. Method by Pixels

In contrast to the models used in [48], [49] proposed a method of disaggregating populations located in spatial units to the level of pixels. The model is based on an iterated regression procedure, using Dempster's expectation-maximization (EM) algorithm [49], which corresponds to a statistical method for the estimation of parameters in situations of absence of data (or incomplete data). The model allows for the association between population and the reflectance of each pixel of the image, improving the performance of the estimates, especially under extreme conditions of population density.

From the same districts (Ballarat and Geelong), [49] used Landsat TM imagery pixels, initially classified as residential or non-residential. The initial population estimate for each pixel, $p_i$, is given by dividing the population of the individual zone by the number of pixels in that zone, so that the population density in each individual zone is homogeneous or constant. In these zones, the expectation-maximization (EM) algorithm is used to construct and re-estimate an iterated regression of population pixels. First, the regression equation in (3) is estimated and the predicted (estimated) values for $p_i$ are adjusted so that the total population of the individual zone after the adjustment is maintained as the known total population of the zone. According to Reis [43–45], the adjusted population of pixel $i$ is given by the sum of the estimated population of the pixel $i$ and the mean of the residues of the census tract to which the pixel belongs. In the next iteration, the previously adjusted $p_i$'s replace the initial population estimates in the dependent variable and the regression equation is estimated again. The new adjusted $p_i$'s are calculated and replace the current $p_i$'s in the next iteration. The iterations continue and the stopping criterion can be defined according to some measure of fit quality such as the determination coefficient ($R^2$) or the mean square of the residuals, for example. The iterations would end when one of these measures had no changes considered relevant from one iteration to another [43–45].

According to Harvey [49], the multicollinearity between the reflectance in the TM bands can cause convergence problems in this iterated regression process. Another logical problem is that of the negative estimates for the populations associated with the pixels, since the linear regression model used has no restrictions. In an attempt to solve this problem, an alternative procedure was used: At each iteration, the negative estimates were turned to zero and adjustments to the estimates of the other pixels were made in order to keep the total population of the census constant. The predictive validity of the model was also tested from the application of the adjusted regression equation in the second image (from the second district, as in Harvey [48]).

The lower relative error median found for the population in spatial units was 14%. The model based on pixels proved to be more robust than the model based on area, especially in areas of extreme population density, although the effectiveness of the model was not the same in these situations: That is, even in the model at the pixel level, there was a tendency of underestimation in the most densely populated areas and overestimation in the less dense areas.

In Brazil, the works that use remote sensing for population estimates are still scarce. Reis [44] and [46], for example, applied the techniques to the estimation of the population by census tracts in Belo Horizonte, Minas Gerais, for the year 1996 with the help of the TM Landsat 5 sensor (bands 1 to 5 and 7). The results showed a relative median error of 30.4% (versus 14% in Harvey [48]), but much higher at the aggregate level (i.e. at the municipal level), with a total relative error of only −0.06%.

### 2.3. Data Structuring

The structuring of the data consisted of the creation of a database with information from the demographic censuses of 2000 and 2010, and Landsat satellite images for the years 2000, 2010, and 2015. From these data, models of both types (census tracts [48] and pixels [49]) were built from data from 2000 and 2010.

### 2.3.1. Estimates Based on 2000 Data

Regression models at the level of census tracts (as in Harvey [48]) and at the level of pixels (as in Harvey [49]) were created from data from 2000, and were used to construct estimates for the three years: 2000 (internal validation), 2010 (external validation) and 2015 (comparison with other post-census estimates, such as the official IBGE projection). To that end, the dasymetric mapping was carried out for the year 2000, by disaggregating the population by census tracts (from the 2000 census) to the pixel level of the Landsat images. This procedure allowed for compatibility, for the same unit of analysis (the pixels), between the population and surface reflectance information, which was necessary for the construction of the regression models, as in [49]. On the other hand, for the construction of models at the census tract level (as in Harvey [48]), the opposite procedure was carried out: The average of the reflectance from the pixels belonging to each census tract was calculated for each census tract in the 2000 census.

Therefore, the two types of models, based on census tracts [48] and pixels [49], were calculated for the three years, which allowed for the evaluation of the results both in the context of internal and external validation. External validation is a procedure not yet performed in the literature among the works that use Harvey's methods [48,49]. The external validation of the estimates calculated by Harvey was carried out for the same year, while the validation step of [33] was restricted to internal validation.

### 2.3.2. Estimates Based on 2010 Data

Regression models were created from data from 2010, also at the level of the census tracts [48] and pixels [49], and used for the construction of estimates for two years: 2010 (internal validation) and 2015. Dasymetric mapping was carried out for the year 2010 by disaggregating population data from two sources: census tracts (2010 census) and the statistical grids cells [31], using the method based on the level of pixels for each [49]. As in the case of the 2000 models, this procedure allowed the compatibility, at the pixel level, of the population and surface reflectance information necessary for the construction of the regression models [49], based on data from 2010. For the construction of the models at the level of census tracts [48], the average reflectance coming from the pixels belonging to each census tract in 2010 was calculated.

It is important to note that there were changes in the administrative boundaries of the census tracts between 2000 and 2010. During the period, new census tracts were created, from the subdivision of pre-existing census tracts in 2000, in addition to changes in the geographical boundaries of a portion of the census tracts. This inconsistency was not a problem, as estimates were calculated for the entire municipality. The changes occurred at the boundary of the census tracts and within the municipality. However, such inconsistency was actually positive because it allowed for the comparison between the models adjusted from data from 2000 and 2010 in relation to the differences in the calculated estimates. One hypothesis is that the models calculated with data from 2010, which are based on a larger number of census tracts, will present better results of internal validation when compared with models calculated with data from 2000. The same logic applies to the models calculated from the 2010 statistical grids cells. These models are expected to have a better fit when compared to the models from the 2000 data.

The estimates calculated for the year 2015—for the models built with data from 2000 and 2000 were compared with four data sources: IBGE projection (by the method of demographic components for the state level, and disaggregated for the municipality by the method AiBi) and by three types of demographic methods that use symptomatic variables to calculate the estimates (simple extrapolation, vital rates and census ratio). Given the absence of census surveys, the 2015 data sources were restricted to estimates. Thus, the comparisons between the 2015 estimates did not constitute "external validation". Nevertheless, the comparison between different methods can contribute to the analysis of effectiveness of estimates produced via orbital images. Very similar results between the different methods, for example, can be analyzed in the light of the costs related to the time and processing required to calculate

the estimates via remote sensing. At the same time, very different results can be indicative of possible problems and/or limitations of some methods.

The simple extrapolation method was chosen for comparison with the remote sensing estimates produced for 2015 due to the extremely short time required for the calculation and under the hypothesis that this method is suitable for the work proposed, which is to produce estimates for short intervals of time (in post-census periods). The vital rate method, in turn, was chosen for comparison due to the advances achieved in recent years in terms of reducing the under-enumeration of birth and death records. Finally, the census ratio method was chosen for comparison with the remote sensing estimates produced for 2015 because it is a method that uses only one symptomatic variable. For this method, this study chose to use information on enrollment in elementary school, obtained from the School Census, motivated by the fact that Brazil achieved the universalization of primary education a few years ago.

### 2.4. Dependent and Explanatory Variables

Table 2 shows the set of variables for the construction of models based on Harvey [48], at the level of the census tracts.

**Table 2.** Variables used for the construction of models based on census tracts.

| Abbreviation | Variables | Source |
|---|---|---|
| Dens00 | Density of the Census tract in 2000 | IBGE |
| Dens10 | Density of the Census tract in 2010 | IBGE |
| Log(Dens)00 | Logarithm of the density of the census tract in 2000 | IBGE |
| Log(Dens)10 | Logarithm of the density of the census tract in 2010 | IBGE |
| TM1 | Average reflectance in the band 1 (Landsat 7 ETM+) | USGS |
| TM2 | Average reflectance in the band 2 (Landsat 7 ETM+) | USGS |
| TM3 | Average reflectance in the band 3 (Landsat 7 ETM+) | USGS |
| TM4 | Average reflectance in the band 4 (Landsat 7 ETM+) | USGS |
| TM5 | Average reflectance in the band 5 (Landsat 7 ETM+) | USGS |
| TM7 | Average reflectance in the band 7 (Landsat 7 ETM+) | USGS |
| AreaOcup00 | Percentage of occupied area in 2000 | USGS/IBGE |
| AreaOcup10 | Percentage of occupied area in 2010 | USGS/IBGE |
| AgSub00 | Subnormal agglomerate in 2000 | IBGE |
| AgSub10 | Subnormal agglomerate in 2010 | IBGE |

Below is the list of the 7 models at the level of the census tracts that were based on [44,46,48]:

- Model 1: dependent variable = density; explanatory variables: bands 1 to 5 and 7;
- Model 2: dependent variable = density logarithm; explanatory variables: bands 1 to 5 and 7;
- Model 3: dependent variable = density logarithm; explanatory variables: bands 1 to 5 and 7 and percentage of occupied area;
- Model 4: dependent variable = density logarithm; explanatory variables: bands 1 to 5 and 7 and percentage of occupied area and subnormal agglomerate (slums);
- Model 5: dependent variable = density logarithm; explanatory variables: bands 1, 4, 5 and 7, percentage of occupied area and subnormal agglomerate (slums);
- Model 6: dependent variable = density logarithm; explanatory variables: bands 1, 4 and 7, percentage of occupied area and subnormal agglomerate (slums);
- Model 7: dependent variable = density logarithm; explanatory variables: bands 1, 4 and 7, percentage of occupied area and subnormal agglomerate (slums) and three interaction variables between bands 1, 4 and 7 and subnormal agglomerate(slums).

First, the models were tested with the incorporation of all bands, as observed in [48]. However, the best adjustments were verified in the models in which the dependent variable was the density logarithm. Thus, models 2 to 7 incorporated this change, with model 3 incorporating the variable

percentage of occupied area (which, in many cases, improved the predictive power of the models, as discussed in the results). In model 4, the subnormal agglomerate (slums) variable was incorporated; in model 5, bands 2 and 3 were removed; in model 6, band 5 was removed; and in model 7, interaction variables between bands 1, 4 and 7 with subnormal agglomerate were incorporated, in order to verify the influence of the slum on the effect that each band has on density.

The linear correlation was high between the reflectances of bands 1, 2 and 3, as well as between bands 5 and 7, while band 4 showed the lowest correlations with the other bands, as commented in the results. According to Reis [44], these relationships between independent variables can cause multicollinearity problems, affecting the estimation of the model coefficients. In general, the removal of bands 2, 3 and 5 resulted in better adjusted models [44].

Table 3 shows the set of variables for the construction of models based on Harvey [49], at the pixel level.

**Table 3.** Variables used for the construction of models based on pixels.

| Abbreviation | Variables | Source |
|:---:|:---:|:---:|
| Pop00pixel(census tract) | Population of the pixel in 2000 | IBGE |
| Pop10pixel(cens tract or grid) | Population of the pixel in 2010 | IBGE |
| TM1 | Reflectance of the pixel in the band 1 (Landsat 7 ETM+) | USGS |
| TM2 | Reflectance of the pixel in the band 2 (Landsat 7 ETM+) | USGS |
| TM3 | Reflectance of the pixel in the band 3 (Landsat 7 ETM+) | USGS |
| TM4 | Reflectance of the pixel in the band 4 (Landsat 7 ETM+) | USGS |
| TM5 | Reflectance of the pixel in the band 5 (Landsat 7 ETM+) | USGS |
| TM7 | Reflectance of the pixel in the band 7 (Landsat 7 ETM+) | USGS |
| Urban00 | Situation (Urban / Rural) in 2000 | IBGE |
| Urban10 | Situation (Urban / Rural) in 2010 | IBGE |
| AgSub00 | Subnormal agglomerate in 2000 | IBGE |
| AgSub10 | Subnormal agglomerate in 2010 | IBGE |

Below is the list of the 6 models at the level of pixels that were based on [44,46,49]:

- Model 1: dependent variable = population; explanatory variables: bands 1 to 5 and 7;
- Model 2: dependent variable = population; explanatory variables: bands 1 to 5 and 7 and subnormal agglomerate (slums);
- Model 3: dependent variable = population; explanatory variables: bands 1 to 5 and 7, subnormal agglomerate (slums) and situation (urban / rural);
- Model 4: dependent variable = population; explanatory variables: bands 1, 4, 5 and 7, subnormal agglomerate (slums) and situation (urban / rural);
- Model 5: dependent variable = population; explanatory variables: bands 1, 4 and 7, subnormal agglomerate (slums) and situation (urban / rural);
- Model 6: dependent variable = population; explanatory variables: bands 1, 4 and 7, subnormal agglomerate (slums), situation (urban / rural) and three interaction variables between bands 1, 4 and 7 and subnormal agglomerate(slums).

Similarly, the models based on [49] started from model 1, which presented all bands as explanatory variables. Model 2 was the incorporation of the variable referring to subnormal agglomerate, which generally better predicts the models' power, while model 3 incorporated the urban/rural variable. Model 4 corresponded to the exclusion of bands 2 and 3, while model 5 referred to the exclusion of band 5. Finally, model 6 incorporated three interaction variables between bands 1, 4 and 7 and subnormal agglomerate.
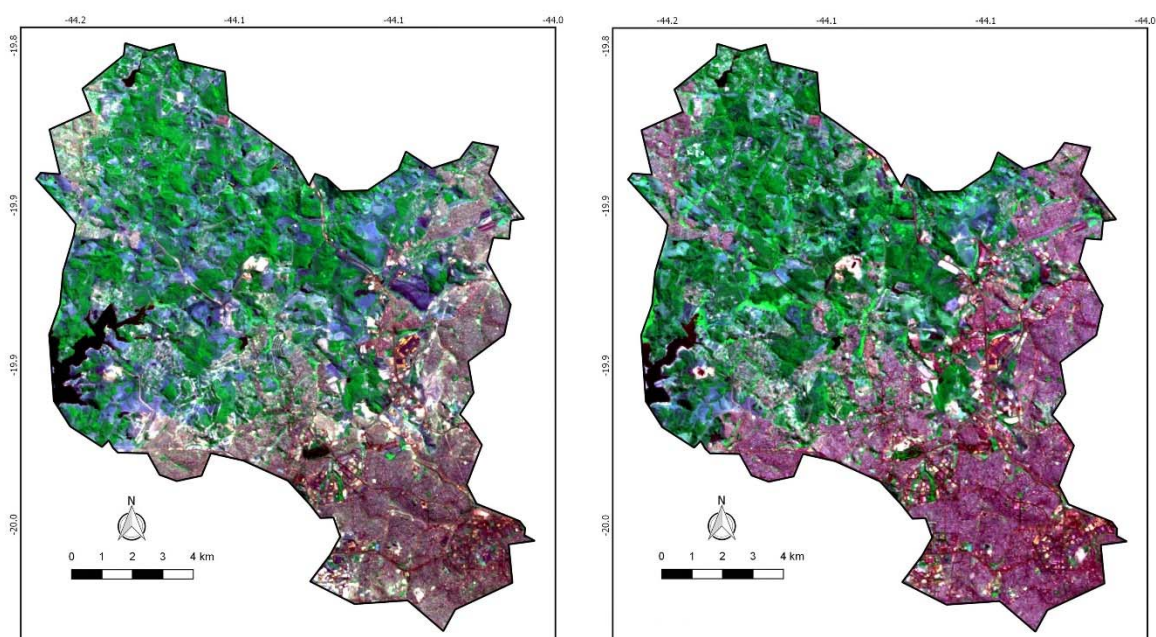
*2.5. Error Measures*

The error measures used for the analysis of the internal validation were the median relative error, $R^2_{back}$ and the total error. The median relative error corresponded to the median of the absolute values of the relative errors, observed for each census tract (in the applications of [48]) and pixel (in the case of [49]). The second measure was $R^2_{back}$, which corresponded to the square of the linear correlation coefficient between population estimates for tracts (or pixels) and the actual population values of those tracts (or pixels). The closer to one (1), the greater the fit of the model. As it is calculated at the sector level, $R^2_{back}$ is more connected to the median relative error than to the total error, and therefore tends to have low values. The third measure was the total relative error, which represented the variation of the estimated total in relation to the total observed for the set of tracts (or pixels), that is, for the municipality as a whole [48,49].
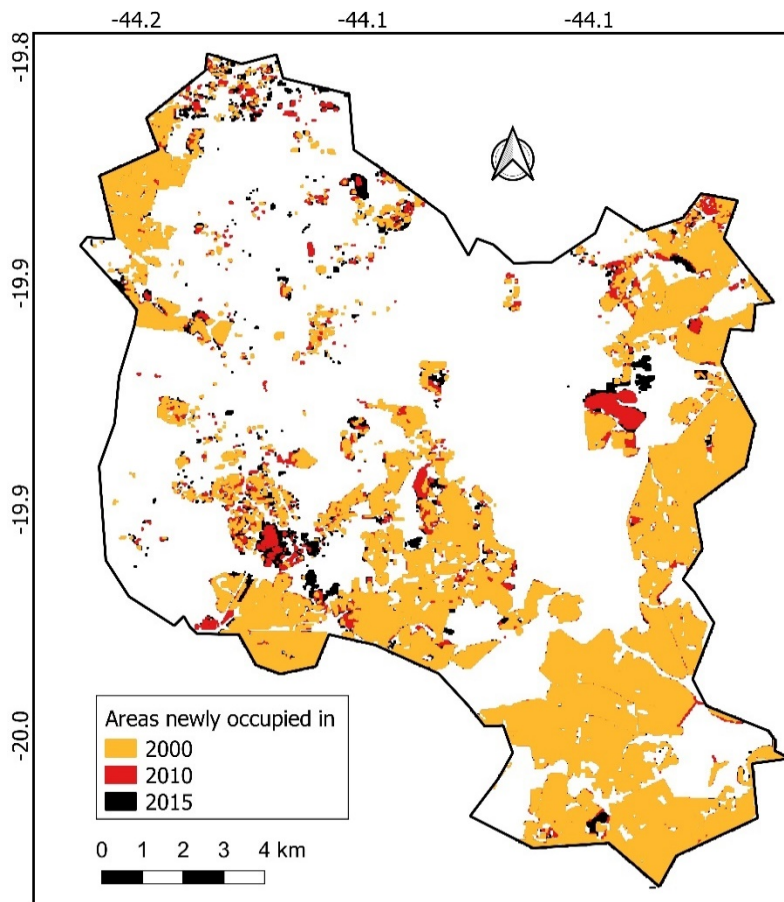
## 3. Results

Figure 1 shows the RGB345 composition of the Landsat ETM+ satellite imagery for Contagem in 2000 and 2015, where the urban concentration in the southeast and eastern regions of the municipality is observed. In the central, southwest and northern regions of the municipality, there is a predominance of areas of low population density, marked by rural areas and absence of occupation. Figure 2 complements Figure 1 showing the evolution of occupancy spots between 2000 and 2015.
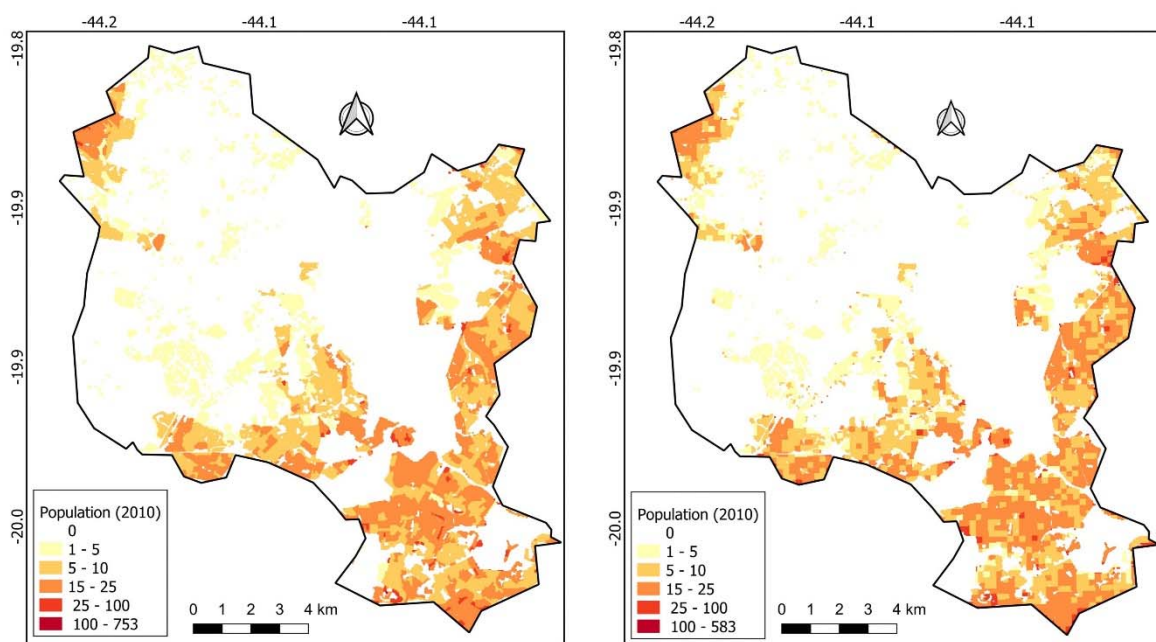
In 2000, 2010 and 2015, 57,174, 63,857 and 69,565 pixels were classified, respectively, as human occupation areas, representing a growth of 11.7% in the period 2000-2010 and 8.9% between 2010 and 2015. Among the total of pixels classified as areas occupied in the period, 95.5% (2000), 96.5% (2010) and 95.2% (2015) are located in urban census tracts. Unlike Reis [44], the rural census tracts were not excluded from the databases for the production of the estimates, which allowed the incorporation of the explanatory variable "urban/rural" in the study. Regarding subnormal agglomerate, only 5% of the pixels were located in slums in 2000, a percentage that increased to 5.6% in 2010 and 5.4% in 2015. Figure 3 shows the dasymetric mappings from the census tracts of 2010 and the 2010 statistical grids cells.



**Figure 1.** RGB345 composition of the Landsat ETM+ satellite imagery in 2000 (left) and 2015 (right) - municipality of Contagem.

**Figure 2.** Evolution of the occupied areas (2000, 2010 and 2015) from the classification of Landsat ETM+ satellite imagery - municipality of Contagem.



**Figure 3.** Dasymetric mapping through Landsat 7 ETM+ satellite imagery pixels. Data from the census tracts of the 2010 census (left) and the statistical grids cells of 2010 (right) – municipality of Contagem.

### 3.1. Models and Estimates at the Level of Census Tracts

For the definition of explanatory variables and models at the level of the census tracts (as in Harvey [48]), this study first verified the correlation between the bands of the selected Landsat ETM+ satellite imagery. A high correlation between the neighboring bands (as expected, due to the overlap between the bands of the spectrum) was observed, especially between bands 1, 2 and 3 (close to 1). In 2010, the behavior of the correlations between density and bands was quite similar, as well as in the models based on pixel (which, for this reason, are not presented in this paper).

For the construction of the models, 2/3 of the databases were randomly defined to fit each model, while the other 1/3 were used for the internal validation test. Table 4 presents the internal validation results of models at the level of census tracts, constructed from the 2000 and 2010 data.

**Table 4.** Internal validation of models at the census tract level (2000 and 2010) - municipality of Contagem.

|  | Indicators | Model 1 | Model 2 | Model 3 | Model 4 | Model 5 | Model 6 | Model 7 |
|---|---|---|---|---|---|---|---|---|
| 2000 (sector) | $R^2$ (back) | 0.061 | 0.032 | 0.232 | 0.323 | 0.322 | 0.341 * | 0.341 |
|  | Median relative error | 0.361 | 0.399 | 0.346 | 0.349 | 0.330 | 0.340 * | 0.321 |
|  | Total error (%) | 115.7 | 44.40 | 10.52 | 6.56 | 6.81 | 6.29 * | 6.18 |
| 2010 (sector) | $R^2$ (back) | 0.262 | 0.179 | 0.281 | 0.322 | 0.335 | 0.341 * | 0.341 |
|  | Median relative error | 0.316 | 0.386 | 0.305 | 0.308 | 0.292 | 0.300 * | 0.298 |
|  | Total error (%) | 25.82 | 17.42 | −0.26 | −1.82 | −1.76 | −1.50 * | −1.52 |

\* Model used to calculate the estimates.

In general, $R^2_{back}$ was relatively weak for all models, either in 2000 or 2010. It is important to remember that, because it is calculated at the sector level, $R^2_{back}$ is more related to the average error on census tracts than to total error. In the total error (macro level), census tracts in which the population was overestimated were compensated for by census tracts in which the population was underestimated. As the total error works with the sum of the estimated populations, the individual errors of the census tracts end up being diluted. $R^2_{back}$ does not have this chance of "error compensation" because the calculation "sees" the individual values of the census tracts. Regarding the median relative error (MRE) of the census tracts, this study observed little difference between models (ranging between 0.321 and 0.399 in 2000, and 0.292 and 0.386 in 2010).

In the internal validation, the models that obtained the best results were model 7 (with a total error of 6.18% in 2000) and model 3 (with a total error of −0.26% in 2010). The use of the logarithm in the dependent variable drastically reduced the total error (from 115.7% to 44.4% in 2000 and from 25.82% to 17.42% in 2010). Similarly, the incorporation of the variable occupied area (%) (model 3) greatly improved the results to levels of approximately 10.52% in 2000 and −0.26% in 2010. The incorporation of the subnormal agglomerate variable decreased the total error from 10.5% to 6.56% in 2000, while in 2010 the error increased from −0.26% to -1.82%. At much lower intensity, the exclusion of bands 2, 3 and 5 and the exclusion of interaction terms (from model 7) did not significantly improve the predictive power of the models.

Table 5 shows the results of the external validation, through the percentage differences between the 2010 remote sensing estimate (constructed using the model created in 2000 at the tract level) and the 2010 census population. The table also shows the percentage differences of the Contagem population, obtained through IBGE projections (revisions 2008 and 2013) with the 2010 census population.

In 2000, although the lowest total error was observed in model 7 (6.18%), model 6 had a very similar total error (6.29%), but without the need to incorporate the interaction variables between the bands and subnormal agglomerate. Therefore, model 6 was used to calculate the estimates for the external validation (2010) and for comparison with other demographic estimates (2015). Regarding 2010, although model 3 presented the lowest total error (−0.26%), model 6 presented a very small total error (−1.5%) and a MRE lower than observed in model 3 (0.300 vs. 0.305). Thus, model 6 was adopted as a reference for calculating the 2010 and 2015 estimates.

The 2000 model, when applied to the 2010 data, resulted in a MRE of 27.0%, while the 2010 estimates, produced from this model, showed a total error of 11.1% (670,287), when compared to the population of Contagem (603,442, obtained in the 2010 Census). This error can be considered acceptable [51,52]. On the other hand, the disaggregation of the IBGE projections (revisions 2008 and 2013) to the municipal level of Contagem (via the application of the AiBi method) presented total errors of 5.0% (633,361) and 4.5% (630,352), respectively; that is, errors lower than that of the estimate produced by remote sensing (calculated by models constructed with data at the level of the 2000 census tracts).

**Table 5.** External validation: Percentage differences of the estimates via remote sensing for the year 2010 (based on a model created in 2000, at the sector level) and the official estimates of 2010, in relation to the population of the 2010 census - municipality of Contagem.

| Municipality | Abs/Diff (%) | Census 2010 | Pop of 2010 (Projection IBGE 2013) | Pop de 2010 (Projection IBGE 2008) | Model of 2000 - Census Tracts |
|---|---|---|---|---|---|
| Contagem | Absolute | 603,442 | 630,352 | 633,361 | 670,287 |
| | Difference (%) | 0.0 | 4.5 | 5.0 | 11.1 |

It is important to emphasize that the good results found in the internal validation stage are in line with international and national references on the subject. On the other hand, the external validation of this study is an innovation, since the external validation performed by [48] is restricted to the application of models for the same year (and another area). In addition, the external validations proposed in this paper were tested for another census year (and, consequently, from another image), which increased the challenge in relation to the assertiveness of the estimates.

The 2000 and 2010 models were used for the construction of estimates for the year 2015, and, for comparative analysis, were used as parameters to post-census IBGE estimates for the municipality of Contagem (which, in 2015, was 648,766 people). Table 6 shows the percentage differences between the estimates by remote sensing calculated through models at the level of census tracts (based on data from 2000 and 2010) and the IBGE post-census estimate; it also compared these estimates (IBGE post-census estimate) with the demographic estimates by simple methods of extrapolation, vital rates and census ratio.

**Table 6.** Comparative analysis: Percentage differences of the estimates via remote sensing (based on models created in 2000 and 2010, at the tract level) and demographic estimates, for the year 2015, in relation to the IBGE's post-census estimates - municipality of Contagem.

| Municipality | Abs/Diff (%) | IBGE Post-Census Estimate (2015) | Simple Extrapolation | Vital Rates | Census Ratio | Model of 2000 - Census Tracts | Model of 2010 - Census Tracts |
|---|---|---|---|---|---|---|---|
| Contagem | Absolute | 648,766 | 636,155 | 666,439 | 622,170 | 701,389 | 691,204 |
| | Difference (%) | 0.0 | −1.9 | 2.7 | −4.1 | 8.1 | 6.5 |

The model based on data from 2000 presented an estimated 701,389 people (difference of 8.1% in relation to the IBGE estimate), while the model based on data from 2010 estimated the population at 691,204 (difference of 6.5% in relation to IBGE's estimate). Demographic estimates by simple extrapolation, vital rates, and census ratios showed percentage differences of −1.9%, 2.7%, and −4.1%, respectively, smaller errors compared to the estimates by remote sensing.

However, the search for assertiveness in relation to the population of 2015 is not a good alternative. It should be considered that the IBGE population for 2015 is only a post-census estimate, calculated by disaggregation of the projection by IBGE components (revised 2013) at the municipal level using the AiBi method, and has limitations (like any extrapolation method). The AiBi method projects the population of a small area from its contribution to the absolute population growth expected in the larger area, and assumes a linear relationship between the growth of the larger area and the smaller area,

which may not be verified, especially in smaller areas, such as municipalities. The idea of comparing post-census estimates in the year 2015 with demographic projections is to analyze only whether there is a large discrepancy with official estimates.

*3.2. Models and Estimates at the Level of Pixels*

As in models based on census tracts, for the construction of the models at the pixel level, 2/3 of the databases were randomly selected. Table 7 shows the internal validation results of models based on pixels, constructed from data from 2000 (by census tracts) and 2010 (by census tracts and statistical grids cells).

In all models based on pixels (data of census tracts 2000 and 2010, and statistical grids cells 2010), the $R^2_{back}$ showed relatively low values. This refers to the weak relation between population and the explanatory variables at the level of the pixels. This is corroborated by high MREs at the pixel level. However, when analyzing the internal validation of the models, the estimates show very low total errors, lower than 2% in all models in the three databases under analysis, which can be considered an excellent result.

**Table 7.** Internal validation of models at the pixel level (2000 and 2010) - municipality of Contagem.

| | Indicators | Model 1 | Model 2 | Model 3 | Model 4 | Model 5 | Model 6 |
|---|---|---|---|---|---|---|---|
| 2000 (sector) | $R^2$ (back) | 0.077 | 0.184 | 0.110 | 0.126 * | 0.110 | 0.146 |
| | Median relative error | 0.353 | 0.315 | 0.317 | 0.314 * | 0.316 | 0.318 |
| | Total error (%) | 1.73 | 0.09 | 1.61 | 0.02 * | 1.67 | 0.84 |
| 2010 (sector) | $R^2$ (back) | 0.144 | 0.186 | 0.154 | 0.171 | 0.170 | 0.269 * |
| | Median relative error | 0.312 | 0.293 | 0.302 | 0.282 | 0.281 | 0.280 * |
| | Total error (%) | 0.15 | 0.68 | 1.85 | 1.28 | 1.17 | −0.40 * |
| 2010 (statistical grid) | $R^2$ (back) | 0.099 | 0.127 | 0.136 | 0.110 | 0.170 * | 0.141 |
| | Median relative error | 0.330 | 0.323 | 0.328 | 0.315 | 0.318 * | 0.319 |
| | Total error (%) | -0.28 | −1.19 | −1.16 | −1.22 | 0.40 * | −1.34 |

* Model used to calculate the estimates.

The models at the level of pixels using data from 2000 census tracts had the highest $R^2$ observed in model 2 (0.184), followed by models 6 (0.146) and 4 (0.126), while the lowest MRE was detected in model 4 (0.314). The lowest total error was also noticed in model 4 (0.02%), below even the results found by [44], of −0.06%. This is considered an excellent result. Therefore, model 4 was chosen for the calculation of estimates from the data of census tract 2000. It should be noted that the incorporation of the variable subnormal agglomerate (model 2) greatly improved the total estimated error, while the incorporation of the variable urban/rural situation (model 3) did not improve the calculated estimate. On the other hand, model 6 was the one that obtained the best results among those that use data from the 2010 census tracts. With $R^2_{back}$ of 0.269, model 6 had the lowest MRE (0.280) and the lowest total error (−0.40%). It is worth noting that the incorporation of the variables subnormal agglomerate (model 2) and urban/rural situation (model 3) did not improve the total error, when compared to the model that used only the reflectance of bands 1 to 5 and 7 (model 1).

The models that used statistical grids cells obtained the best fit observed in model 5, with 0.170 of $R^2_{back}$, 0.318 of MRE, and 0.4% of total error (being chosen to calculate the estimate). Neither the incorporation of the variables subnormal agglomerate (model 2) and urban/rural situation (model 3) nor the use of data from the census tracts of 2010 improved the total error, when compared to the model that used only the reflectance of bands 1 to 5 and 7 (model 1). The results show that the models based on pixels fit better when compared to models at the level of the census tracts, especially in relation to the internal validation. Total errors were found to be less than 0.5% in the models at the level of pixels, against 6.29% and -1.5% in the models at the level of the census tracts for the 2000 and 2010 databases.

In addition, the indicators showed that the database by statistical grids cells did not present the best results, although the difference was very small. The results were similar to the models constructed by 2010 census tracts. It is also observed that, in the three databases under analysis, the models with

the best internal validations (4, 5 and 6) are those that do not consider bands 2 and 3, as well as those that use iteration terms variables.

The model 4 at the level of pixels, based on data from the census tracts of 2000, was used to produce estimates for the year 2010, in order to perform the external validation (Table 8). On the other hand, model 6, based on data from the census tracts of 2010, and model 5, based on the statistical grids cells of 2010, were used to produce estimates for the year 2015 (Table 9). As already mentioned, the objective is to carry out a comparative analysis of these post-census estimates based on pixels with demographic estimates, having as reference the post-census estimates of the IBGE.

**Table 8.** External validation: Percentage differences of the estimates via remote sensing for the year 2010 (based on a model created in 2000, at the pixel level) and the official estimates of 2010, in relation to the population of the 2010 census - municipality of Contagem.

| Municipality | Abs/Diff (%) | Census 2010 | Pop of 2010 (Projection IBGE 2013) | Pop of 2010 (Projection IBGE 2008) | Model of 2000 - Pixels |
|---|---|---|---|---|---|
| Contagem | Absolute | 603,442 | 630,352 | 633,361 | 599,371 |
| | Difference (%) | 0.0 | 4.5 | 5.0 | −0.7 |

Regarding the external validation, the application of the model at the pixel level based on data from 2000 to the production of estimation in 2010 resulted in MRE of the census tracts of 28.9%, and an estimated population of 599,371, which represented a percentage difference of only −0.7% in relation to the population of Contagem. The IBGE projections (revision 2008 and 2013), disaggregated by the AiBi method for the municipal level, presented a total error of 5.0% (633,361) and 4.5% (630,352), respectively. This shows, at first, the best fit of the estimates calculated at the pixel level, via Landsat ETM+ satellite imagery and data from the 2000 Census tracts.

With regard to the comparative analysis for the year 2015, Table 9 shows the results of the models based on pixels and the comparisons with the IBGE estimate and the post-census demographic estimates.

**Table 9.** Comparative analysis: Percentage differences of the estimates via remote sensing (based on models created in 2000 and 2010, at the pixel level) and demographic estimates, for the year 2015, in relation to the IBGE's post-census estimates - municipality of Contagem.

| Municipality | Abs/Diff (%) | IBGE Post-Census Estimate (2015) | Simple Extrapolation | Vital Rates | Census Ratio | Model of 2000 - Pixels (From Census) | Model of 2010 - Pixels (From Census) | Model of 2010 - Pixels (From Statistical Grids) |
|---|---|---|---|---|---|---|---|---|
| Contagem | Absolute | 648,766 | 636,155 | 666,439 | 622,170 | 701,199 | 664,205 | 683,835 |
| | Difference (%) | 0.0 | −1.9 | 2.7 | −4.1 | 8.1 | 2.4 | 5.4 |

The estimates using orbital images Landsat 7 for the year 2015 resulted in estimated populations and total errors of 701,199 and 8.1% (2000 model, from census tracts), 664,205 and 2.4% (2010 model, from census tracts), and 683,835 and 5.4% (2010 model, from statistical grids cells), which represented excellent results, although the errors were larger when compared to the demographic estimates.

## 4. Discussion

### 4.1. About Internal Validation

For the municipality of Contagem, marked by high population density, the results showed the potential of the methods used for various purposes in the field of estimation for small areas. First, very low errors were found in the internal validation step for both models at the level of census tracts (6.29% for the 2000 model and -1.5% for 2010 model), and for all models at the pixel level (0.02% for 2000 model, −0.4% for the 2010 model from census tracts, and 0.4% for the 2010 statistical grids cells model). These results are in line with the literature. The studies of [44,46] presented a relative

error in the pixel method of −0.06%, while in [48] the total relative error was -4.8%. The results show the potential of these methods (especially the one based on [49], at the pixel level) for estimating the distribution of the population in the continuous space and in the production of intercensal estimates, via Landsat imagery. The good results of the internal validation for Contagem also demonstrate the great potential to estimate areas of high density for the same year and the same scene, which can be very useful in areas where census coverage has not been satisfactory.

*4.2. About External Validation*

This paper contributes to the existing literature on carrying out external validations between censuses through the construction of models in one year and the calculation and analysis of error of estimates in another year, over a period of 10 years. The results obtained in steps of external validation and comparative analysis, especially in models of the pixels, show the potential of this method for the production of reliable post-census estimates for cities with high population density.

The 2000 model, when applied to the 2010 data, showed a total error of 11.1% (670,287), when compared to the population of Contagem, obtained in the 2010 Census. Therefore, the model at the pixel level presented external validation of −0.7%, while in the comparative analysis with the IBGE post-census estimate (in 2015), the total errors were between 8.1% and 2.4% (respectively, to the 2000 and 2010 models by census tracts). This may be considered acceptable, although they were at levels higher than the demographic estimates, which were between 2% and -4%. It should be remembered that the real population of 2015 is not known, but the small percentage differences represent a good indication of the adjustment of the models by orbital images.

The external validation obtained through a model at the pixel level, based on data from 2000 (−0.7%), does not present a pattern of similarity to the total error of the same model calculated for the year 2015 (8.1%), which may be associated with either the error of the IBGE post-census estimates for 2015 or, more likely, the difficulty of obtaining an error pattern between images of different years applied to the same model. This may have occurred because of the different responses that reflectance can provide for different images at different time periods, or because of the two factors together.

**5. Conclusions**

Our results show that there is no single model that can be used (census tracts and pixels), or application (internal and external validation), or for each image. In fact, the quality of the estimates at the level of pixels is conditioned on the maintenance of the relationship between population and reflectance (of spectral bands) over time, between different images, and the pattern of change in this relation seems to be not well established. The changes in the relationship between reflectance and population can occur due to the differences in the level of solar illumination, as well as due to atmospheric interference (such as the differences in the level of cloud cover, although images with low cloud percentage were used), or even due to the loss of 22% of information in Landsat 7 images from 2003 (referring to a noise problem), which, to some extent, can compromise the results.

Regarding the limitations, it was possible to verify that there is no model (or, in other words, a set of explanatory variables) that can be defined as the standard to be replicated. The best fit varied according to the method used, but also according to the image with which the model was constructed. Some limitations of the results found may be associated with the choices defined in this paper. The inefficiency of the database by statistical grids cells in the production of better-quality estimates, when compared to the databases by census tracts, serves as an agenda for future studies. The way forward is to focus processing efforts on a larger number of municipalities to increase the variability of the results in order to contribute to more assertive conclusions, for example, on the possibility of replicability of the method to larger areas.

Another limitation concerns the scale at which estimates are produced. If, on the one hand, dasymetric mapping allows for the calculation of estimates for any administrative unit (including intra-municipal), the MREs of the sectors—with results of approximately 35%, close to that found by [48], of 30.4% and above

that found in [49] of 14%—prevent estimates being produced for highly disaggregated units. Other studies have used Landsat imagery to calculate population estimates [13,14,33,43,53,54]. However, in these studies, different methodologies were used, either in relation to the classification of the area occupied, or the statistical model adopted, which makes it impossible to directly compare the results of those studies with the results of this study for an analysis of the limitations and contributions of the proposed methodology.

However, in spite of the limitations, either from the point of view of the estimates themselves, as a tool for sustainable economic, social, and environmental planning, or as an auxiliary tool in data analysis, for the definition of trends in population projection studies, the results for the municipality of Contagem can be considered promising. These methods have a range of application possibilities for areas of high population density, and great potential for the estimation of small areas in several fields of knowledge.

**Author Contributions:** Conceptualization, J.C., J.I.R.R., E.A.B., A.M.V.M. and I.A.R.; Formal analysis, J.C. and E.A.B.; Funding acquisition, J.C., J.I.R.R. and E.A.B.; Investigation, J.C.; Methodology, J.C., J.I.R.R., A.M.V.M. and I.A.R.; Project administration, J.C.; Resources, J.I.R.R. and E.A.B.; Software, J.C.; Supervision, J.C., J.I.R.R., A.M.V.M. and I.A.R.; Validation, J.C., J.I.R.R., E.A.B., A.M.V.M. and I.A.R.; Visualization, J.C., J.I.R.R., E.A.B., A.M.V.M. and I.A.R.; Writing – original draft, J.C.; Writing – review & editing, J.C. and E.A.B. All authors have read and agreed to the published version of the manuscript.

**Conflicts of Interest:** On behalf of all authors, the corresponding author states that there is no conflict of interest.

## References

1. Shryock, H.S.; Siegel, J.S.; Larmon, E.A. *The Methods and Materials of Demography*, 2nd ed.; Population Estimates; Government Printing Office: Washington, DC, USA, 1973; pp. 725–809.

2. Bryan, T. Population estimates. In *The Methods and Materials of Demography*; Siegel, J.S., Swanson, D.A., Eds.; Elsevier Academic Press: San Diego, CA, USA, 2004; pp. 523–560.

3. Smith, S.K.; Morrison, P.A. Small-area and business demography. In *Handbook of Population*; Poston, D., Micklin, M., Eds.; Springer Publishers: New York, NY, USA, 2005.

4. Cabrera, M. Estimación de Población em áreas menores con métodos que utilizan variables sintomáticas. *Urug. Com. Sect. Poblacíon* **2011**, *61*, 23–27.

5. Jannuzzi, P.M. Cenários futuros e projeções populacionais para pequenas áreas: Método e aplicação para distritos paulistanos 2000–2010. *Rev. Bras. Estud. Popul.* **2007**, *24*, 109–136. [CrossRef]

6. Silva, L.G.C.; Santos, R.O. A utilização de variáveis sintomáticas para estimativas de populações municipais do estado do paraná. *Ann. Encontro Nac. Estud. Popul.* **2016**, *20*.

7. Esquivel, E.A.C. Variables sintomáticas em las estimaciones poblacionales a nivel cantonal em Costa Rica. *Notas Poblacíon* **2001**, *71*, 51–72.

8. Iiasaka, J.; Hegedus, E. Population estimation from landsat imagery. *Remote Sens. Environ.* **1982**, *12*, 259–272. [CrossRef]

9. Sutton, P.; Roberts, D.; Elvidge, C.; Meij, H.A. Comparison of nighttime satellite imagery and population density for the continental United States. *Photogramm. Eng. Remote Sens.* **1997**, *63*, 1303–1313.

10. Wu, S.; Qiu, X.; Wang, L. Population estimation methods in GIS and remote sensing: A review. *GIScience Remote Sens.* **2005**, *42*, 80–96. [CrossRef]

11. Lu, D.; Weng, Q.; Li, G. Residential population estimation using a remote sensing derived impervious surface approach. *Int. J. Remote Sens.* **2006**, *27*, 3553–3570. [CrossRef]

12. Bhaduri, B.; Bright, E.; Coleman, P.; Urban, M.L. LandScan USA: A high-resolution geospatial and temporal modeling approach for population distribution and dynamics. *GeoJournal* **2007**, *69*, 103–117. [CrossRef]

13. Wu, C.; Murray, A. Population estimation using landsat enhanced thematic mapper imagery. *Geogr. Anal.* **2007**, *39*, 26–43. [CrossRef]

14. Zandbergen, P.A.; Ignizio, D.A. Comparison of dasymetric mapping techniques for small-area population estimates. *Cartogr. Geogr. Inf. Sci.* **2010**, *37*, 199–214. [CrossRef]

15. Mennis, J. Dasymetric mapping for estimating population in small areas. *Geogr. Compass* **2009**, *3*, 727–745. [CrossRef]

16. Dória, V.E.M.; Kampel, S.A.; Monteiro, A.M.V. Estimativa e distribuição espacial da população urbana com imagens de satélite de luzes noturnas: Um estudo para a Região Metropolitana de São Paulo, Brasil, com o sensor Visible/Infrared Imaging Radiometer Suite (VIIRS). *Geografia* **2016**, *41*, 527–548.

17. Silva, F.B.; Poelman, H.; Martens, V.; Lavalle, C. *Population Estimation for the Urban Atlas Polygons*; Joint Research Centre of the European Comission: Ispura, Italy, 2013; 22p.

18. CIESIN. *Gridded Population of the World, Version 4 (GPWv4): Population Count, Revision 10*; NASA Socioeconomic Data and Applications Center (SEDAC): Palisades, NY, USA, 2017.

19. Mckee, J.J.; Rose, A.N.; Bright, E.A.; Huynh, T.; Bhaduri, B.L. Locally adaptive, spatially explicit projection of US population for 2030 and 2050. *Proc. Natl. Acad. Sci. USA* **2015**, *112*, 1344–1349. [CrossRef]

20. Jensen, J.R. *Sensoriamento Remoto do Ambiente: Uma Perspectiva em Recursos Terrestres*; Parêntese: São José dos Campos, Brasil, 2009; p. 598.

21. Souza, I.M. *Análise do Espaço Intra-Urbano para Estimativa Populacional Intercensitária Utilizando Dados Orbitais de Alta Resolução Espacial*; Tesis of Remote Sensing; Instituto Nacional de Pesquisas Espaciais (INPE): São José dos Campos, Brazil, 2004; p. 92.

22. Tomás, L.R. *Inferência Populacional Urbana Baseada No Volume de Edificações Residenciais Usando Imagens IKONOS-II e Dados LIDAR*; Tesisin Remote Sensing; Instituto Nacional de Pesquisas Espaciais (INPE): São José dos Campos, Brazil, 2010; p. 108.

23. Xie, Y.; Weng, A.; Weng, Q. Population estimation of urban residential communities using remotely sensed morphologic data. *IEEE Geosci. Remote Sens. Lett.* **2015**, *12*, 1111–1115.

24. Tobler, W.R. Satellite confirmation of settlement size coefficients. *Area* **1969**, *1*, 30–34.

25. Lo, C.P.; Welch, R. Chinese urban population estimates. *Ann. Assoc. Am. Geogr.* **1977**, *67*, 246–253. [CrossRef]

26. Forest, C. *Estimativas Populacionais e de Crescimento de Áreas Urbanas no Estado de São Paulo, com Utilização de Imagens LANDSAT*; Tesis of Remote Sensing; Instituto Nacional de Pesquisas Espaciais: São José dos Campos, Brazil, 1978; p. 124.

27. Holt, J.B.; Lo, C.P.; Hodler, T.W. Dasymetric estimation of population density and areal interpolation of census data. *Cartogr. Geogr. Inf. Sci.* **2004**, *31*, 103–121. [CrossRef]

28. Eicher, C.L.; Brewer, C.A. Dasymetric mapping and areal interpolation: Implementation and evaluation. *Cartogr. Geogr. Inf. Sci.* **2001**, *28*, 125–138. [CrossRef]

29. Mennis, J. Generating surface models of population using dasymetric mapping. *Prof. Geogr.* **2003**, *55*, 31–42.

30. Mennis, J.; Hultgren, T. Intelligent dasymetric mapping and its application to areal interpolation. *Cartogr. Geogr. Inf. Sci.* **2006**, *33*, 179–194. [CrossRef]

31. Bueno, M.C.D. *Grade Estatística: Uma Abordagem Para Ampliar o Potencial Analítico de Dados Censitários*; Tesis of Demography; Instituto de Filosofia e Ciências Humanas, Universidade Estadual de Campinas: Campinas, Brazil, 2014; p. 263.

32. Kaimaris, D.; Patias, P. Identification and area measurement of the built-up area with the Built-up Index (BUI). *Int. J. Adv. Remote Sens. GIS* **2016**, *5*, 1844–1858. [CrossRef]

33. Langford, M. An evaluation of small area population estimation techniques using open access ancillary data. *Geogr. Anal.* **2013**, *45*, 324–344. [CrossRef]

34. Pirowski, T.; Pomietlowska, J. Distribution of Krakow's population by dasymetric modeling method using urban atlas and publicy available statistical data. *Geomat. Environ. Eng.* **2017**, *11*, 83–95. [CrossRef]

35. European Environment Agency. *Towards an Urban Atlas: Assessment of Spatial Data on 25 European Cities and Urban Areas*; European Environment Agency: Copenhagen, Denmark, 2002; p. 152.

36. Li, X.; Zhou, W. Dasymetric mapping of urban population in China based on radiance corrected DMPS-OLS nighttime light and land cover data. *Sci. Total Environ.* **2018**, *643*, 1248–1256. [CrossRef] [PubMed]

37. Lung, T.; Lubker, T.; Ngochoch, J.K.; Schaab, G. Human population distribution modelling at regional level using very high resolution satellite imagery. *Appl. Geogr.* **2013**, *41*, 36–45. [CrossRef]

38. Nagle, N.N.; Buttenfield, B.P.; Leyk, S.; Spielman, S. Dasymetric modeling and uncertainty. *Ann. Assoc. Am. Geogr.* **2014**, *104*, 80–95. [CrossRef]

39. Mennis, J. Increasing the accuracy of urban population analysis with dasymetric mapping. *City* **2015**, *17*, 115–126.

40. Lloyd, C.T.; Chamberlain, H.; Kerr, D.; Yetman, G.; Pistolesi, L.; Stevens, F.R.; Gaughan, A.E.; Nieves, J.J.; Hornby, G.; Macmanus, K.; et al. Global spatio-temporally harmonised datasets for producing high-resolution gridded population distribution datasets. *Big Earth Data* **2019**, *3*, 108–139. [CrossRef]

41. United States Geological Survey (USGS). Landsat Enhanced Thematic Mapper Plus (ETM+). United States. 2017. Available online: https://lta.cr.usgs.gov/LETMP (accessed on 30 January 2017).

42. Karume, K.; Schmidt, C.; Kundert, K.; Bagula, M.E.; Safina, B.F.; Schomacker, R.; Ganza, D.; Azanga, O.; Nfundiko, C.; Karume, N.; et al. Use of remote sensing for population number determination. *Open Access J. Sci. Technol.* **2017**, *5*, 1–9. [CrossRef]

43. Morton, T.A.; Yuan, F. Analysis of population dynamics using satellite remote sensing and US census data. *Geocarto Int.* **2009**, *24*, 143–163. [CrossRef]

44. Reis, I.A. Estimação da população dos setores censitários de Belo Horizonte usando imagens de satélite. *Ann. Simpósio Bras. Sens. Remoto* **2005**, *12*, 2741–2748.

45. Durand, C.; Pereira, M.N.; Moreira, J.C.; Freitas, C.C. Análise da correlação entre população e área urbana (km$^2$) visando a inferência populacional por meio do uso de imagens orbitais. *Geografia* **2007**, *16*, 113–142.

46. Reis, I.A. Adjusting population estimates using satellite imagery and regression models. *Ann. Simpósio Bras. Sens. Remoto* **2011**, *15*, 830–837.

47. USGS-EROS. 2019. Available online: https://github.com/USGS-EROS/espa-surface-reflectance (accessed on 30 March 2020).

48. Harvey, J.T. Estimating census district populations from satellite imagery: Some approaches and limitations. *Int. J. Remote Sens.* **2002**, *23*, 2071–2095. [CrossRef]

49. Harvey, J.T. Population estimation models based on individual TM pixels. *Photogramm. Eng. Remote Sens.* **2002**, *68*, 1181–1192.

50. Stevens, F.R.; Gaughan, A.E.; Linard, C.; Tatem, A.J. Disaggregating census data for population mapping using random forests with remotely-sensed and ancillary data. *PLoS ONE* **2015**, *10*, 1–22. [CrossRef]

51. Smith, S.K.; Tayman, J.; Swanson, D.A. *State and Local Population Projections: Methodology and Analysis*; Kluwer Academic Publishers: New York, NY, USA; Boston, MA, USA; Dordrecht, The Netherlands; London, UK; Moscow, Russia, 2002; p. 426.

52. Wilson, T.; Rowe, F. The forecast accuracy of local government area population projections: A case study of Queensland. *Australas. J. Reg. Stud.* **2011**, *17*, 204.

53. Li, G.; Weng, Q. Using Landsat ETM+ imagery to measure population density in Indianapolis, Indiana, USA. *Photogramm. Eng. Remote Sens.* **2005**, *71*, 947–958. [CrossRef]

54. Zhang, J.; Xu, W.; Qin, L.; Tian, Y. Spatial distribution estimates of the urban population using DSM and DEM data in China. *ISPRS Int. J. Geo-Inf.* **2018**, *7*, 435. [CrossRef]