A high resolution population grid for the conterminous United States: The 2010 edition

Anna Dmowska^a, Tomasz F. Stepinski^a

^aSpace Informatics Lab, Department of Geography, University of Cincinnati, Cincinnati, USA, OH 45221-0131, USA

Abstract

Readily available high resolution data on population distribution is an important resource for monitoring humanenvironment interactions and for supporting planning and management decisions. Using a grid that approximates population density over the entire country seems like the most practical approach to exploring and distributing detailed population data but instead data based on census aggregation units is still the most widely used method. In this paper we describe the construction of 30 m resolution grid representing the distribution of population in 2010 over the entire conterminous United States. The grid is computed using 2010 U.S. Census block level population counts disaggregated by a dasymetric model that uses land cover (2011 NLCD) and land use (2010 NLUD) as ancillary data. Detailed descriptions of the ancillary data and dasymetric model are given. Methods of computing the grid are presented followed by an extensive assessment of model accuracy. Overall the expected value for relative error of the model is 44% which is at the lower limit of errors reported for other continental-sized, high resolution population grids. We also offer a more specific error estimate for areas with specified value of population density. Using two example areas, one highly urbanized and another rural, we demonstrate the advantages of using the gridded population data over the census block-based data. Our 30 m population grid is available for online exploration and for download from the custom-made GeoWeb application SocScape at http://sil.uc.edu.

19

20

21

22

23

24

25

27

28

29

30

31

32

33

35

36

37

Keywords:

dasymetric modeling, gridded population data, Census, NLCD, land use

1 1. Introduction

Accurate information about the human population distribution is essential for formulating informed re-3 sponses to population-related social, economic, and 4 environmental problems. Governments need precise 5 population data to support planning for infrastructure projects (Benn, 1995; Murray et al., 1998; Pattnaik 7 et al., 1998), locating public facilities (Deng et al., 8 2010), allocating and managing of resources (Gleick, 9 1996; Smith et al., 2002), and for preparing responses 10 to natural disasters (Dobson et al., 2000; Balcik and 11 Beamon, 2008; Maantay and Maroko, 2009; Tenerelli 12 et al., 2015). Similarly, the private sector needs pop-13 ulation data for planning the locations of their facilities 14 (Martin and Williams, 1992), for optimization of service 15 delivery systems, and for risk assessments (Chen et al., 16 2004; Thieken et al., 2006). Reliable information about 17

the population distribution is also essential to assess human pressure on the environment (Weber and Christophersen, 2002), for quantifying environmental impact on population (Vinkx and Visee, 2008), and for public health applications (Hay et al., 2005).

An authoritative source of population data is the government instigated national census; in the U.S. population data is collected every 10 years by the U.S. Census Bureau (hereafter referred to as the census). The most recent census was performed in 2010. The census collects population data with the ultimate resolution of an individual household but it releases this data aggregated to fixed areal units due to privacy concerns. The smallest aggregated areal unit released by the census is the census block. Census blocks in urban areas may be as small as a city block, but they are much larger in suburban and rural areas. There are several reasons why population data aggregated to fixed administrative units is not an ideal form of information about population density.

^{*}Corresponding author. Email address: stepintz@uc.edu

Preprint submitted to Computers, Environment and Urban Systems

First, it suffers from the modifiable areal unit problem 38 (Lloyd, 2014). Second, the spatial detail of aggregated 39 data is variable and low, except in the most densely pop-40 ulated urban areas. Third, there is a spatial mismatch 41 (Voss et al., 1999) between census areal units (blocks, 42 tracts etc.) and user-desired units (for example, neigh-43 borhoods, tax zones, postal delivery zones, vegetation zones, watersheds, etc.). Finally, the boundaries of cen-45 sus aggregation units (particularly blocks) may change 46 from one census to another, making the analysis of pop-47 ulation change at high spatial resolution difficult (Holt 100 48 et al., 2004; Schroeder, 2007; Ruther et al., 2015). 49

Overall, the properties of aggregation unit-based data 102 50 make it ill-suited for the spatial analysis of population- 103 51 related socio-economic and environmental problems. 52 Instead, the population grid has emerged as an alter-53 native format to deliver population data. A population 54 grid is a geographically referenced lattice of square cells 55 with each cell carrying a population count or the value 56 of population density at its location. Population grids 57 are constructed from census unit-based data using ei-58 ther areal weighting interpolation (Goodchild and Lam, 111 59 1980; Flowerdew and Green, 1992; Goodchild et al., 60 1993) or dasymetric modeling (Wright, 1936; Langford 113 61 and Unwin, 1994; Eicher and Brewer, 2001). Popula- 114 62 tion grids have the following advantages: all cells have 115 63 the same size, the cells are stable in time, there is no 116 64 spatial mismatch problem as any partition of the study 117 65 area can be rasterized to be co-registered with a popu-118 66 lation grid. In addition, if dasymetric modeling is used 119 67 (see below), a population grid offers a spatial resolution 120 68 superior to that of the unit-based data. 69

With respect to the construction of a population grid, 122 70 dasymetric modeling can be described as a technique of 123 71 disaggregating aggregation unit-based population data 124 72 into grid cells of a higher spatial resolution using an-125 73 cillary data that correlates with population density but 126 74 which has a higher resolution. Sharpening population 127 75 data using dasymetric modeling has been extensively 128 76 studied (Petrov, 2012) with a focus on the utilization 129 77 of different types of ancillary data in order to increase 130 78 the accuracy of a model. The original, and still the 131 79 most widely used, ancillary data are land cover/land 132 80 use data (Wright, 1936; Mennis, 2003, 2009; Linard 133 81 et al., 2011). High-resolution satellite images have re- 134 82 cently been utilized as ancillary data to identify individ-83 ual buildings (Ural et al., 2011; Lu et al., 2010; Lung 136 84 et al., 2013). A regression analysis is able to link the 85 137 86 area or volume of each building to the number of peo-138 ple in it. If available, the Light Detection and Rang-139 87 ing (LiDAR) data is used (Lu et al., 2010), to help es-140 88 tablish the volume of a building. Another approach to 141 89

dasymetric modeling is to use local infrastructure information, such as street density (Reibel and Bufalino, 2005; Su et al., 2010) or the density of points of interest (Bakillah2014) as ancillary data. Tax parcel data have also been used (Maantay et al., 2007; Kar and Hodgson, 2012; Mitsova et al., 2012; Jia et al., 2014; Jia and Gaughan, 2016) to disaggregate census population data. Other proposed sources of ancillary data include light emission data (Briggs et al., 2007; Sridharan and Qiu, 2013) and address datasets (Zandbergen, 2011).

Despite rapid progress in developing various techniques for dasymetric modeling, the practical adoption of population grids is low. This is because the majority of potential users are only able to utilize the ready-touse product (a population grid) rather then actually create their own. In order to increase the adoption of demographic data for spatial analysis, high resolution grids over broad geographical areas need to be available in the public domain. Such grids have been developed and made available for all countries in the European Union (Gallego, 2010; Gallego et al., 2011) and, through the WorldPop project (http://www.worldpop.org.uk), for countries in South and Central America, Asia and Africa (Gaughan et al., 2013; Linard et al., 2012; Sorichetta et al., 2015). For the United States, the Socioeconomic Data and Application Center (SEDAC) (http://sedac.ciesin.columbia.edu/) provides 1 km resolution (250 m for selected metropolitan areas) demographic grids. However, in addition to having a rather coarse resolution, these grids are only available for the years 1990 and 2000. A higher resolution (90 m) USwide demographic grid, presumably based on most recent census data, is under development by the Oak Ridge National Laboratory (Bhaduri et al., 2007). This project, called LandScan-USA, aims at providing both nighttime (residential) as well as daytime population densities, but it is not currently available, nor is it expected to be in the public domain once it becomes available.

Since 2014 we have been developing high resolution demographic grids for the entire conterminous US. Our goal is to develop grids that offer a significant improvement over SEDAC grids and make them available for exploration and download through our interactive web-based application SocScape (Social Landscape) at http://sil.uc.edu. The first generation of our grids (referred to as SocScape-90) were the results of sharpening SEDAC grids to 90 m resolution using dasymetric modeling with the National Land Cover Dataset (NLCD) as ancillary data (Dmowska and Stepinski, 2014). Using this approach we have developed and made available through SocScape the population grids

an

91

92

93

94

95

96

97

98

99

101

104

107

109

110

for 1990 and 2000. However, our original approach 190 142 had several shortcomings and limitations. First, it did 191 143 not use original census data, instead it relied on the 192 144 SEDAC grid, which, in addition to containing a number 193 145 of errors and inconsistencies (Dmowska and Stepinski, 146 194 2014), was also spatially coarser than census blocks in 147 195 densely populated urban areas. Second, it was limited to 14 years 1990 and 2000 – the only years for which SEDAC 197 149 published its grids. 198 150

In this paper we report on our second generation of 199 151 U.S.-wide grids (referred to as SocScape-30). Our new 200 152 approach differs from the previous approach in the fol-201 153 lowing ways: (1) It uses dasymetric modeling to disag-202 154 gregate census blocks directly, rather then disaggregat-203 155 ing SEDAC cells. (2) It uses two ancillary datasets, the 156 NLCD 2011 and the newly available National Land Use 157 Dataset (NLUD2010) (Theobald, 2014). (3) The new 206 158 grid has a nominal resolution of 30 m, equal to the reso-159 207 lution of both ancillary datasets. (4) We offer an assess-160 ment of uncertainty of the model in the form which is 209 161 directly relevant to a user. SocScape-30 is calculated on 210 162 the basis of 2010 Census block-level data and is avail- 211 163 able online through our GeoWeb application. Section 164 2 describe the datasets used for the construction of the 213 165 2010 grid and section 3 describes our methodology for 214 166 obtaining the population grid. Section 4 gives the de- 215 167 tails of our calculations, presents a quality assessment, 168 and describes how to access the data. Section 5 uses two ²¹⁶ 169 examples, one urban and one rural, to demonstrate the 217 170 advantages of using griided data over the census block-218 171 based data. Discussion and conclusions are given in sec-219 172 tion 6. 173 220

174 2. Input data

The SocScape–30 population grid is constructed using dasymetric modeling. In the context of this paper the dasymetric modeling technique requires two types of data - areal unit-based population data, to be disaggregated to a high resolution grid, and ancillary data at the resolution of this grid.

181 2.1. Census data

The primary source of spatio-demographic infor- 233 182 mation in the United States are decennial censuses 234 183 (http://www.census.gov). The U.S. Census Bureau pro-235 184 vides data as a series of summary text files labeled 236 185 186 from 1 to 4 which provide information at different lev- 237 els of spatial aggregation (from as small as a census 238 187 block to as large as the entire U.S.). To construct 239 188 SocScape-30 we used the population count for each 240 189

block based on variable P1 (total population) from Summary File 1 (SF1). We used the census data distribution provided by the National Historical Geographic Information System (NHGIS) at https://www.nhgis.org/ as we deemed NHGIS distribution easier to use than the Census Bureau distribution. This is because NHGISdistributed demographic tables and shapefiles depicting block boundaries contain identifiers expediting the task of joining population counts to shapefiles.

Sizes of shapefiles containing block boundaries and their population counts vary from 34 MB for the District of Columbia to 4037 MB for the state of California. The overall size of the block-level shapefile/population count data for the entire conterminous U.S. (11007989 blocks) was 39 GB. We converted the block shapefile into a 30 m resolution raster grid co-registered with the ancillary grid (see the next subsection). Grid cells constituting a given block store numerical identifier of this block. There are 8,651,157,015 cells in the grid. Block rasterization is performed in order to expedite the computation of the dasymetric model. An unintended consequence of rasterization is the "loss" of 264565 blocks having a size about equal or smaller than a grid cell. This constitutes about 2% of all the blocks, however, it constitutes a "loss" of only 0.035% of the population because many of these blocks are uninhabited.

2.2. Ancillary data

The role of ancillary data is twofold, first to distinguish between inhabited and uninhabited sections of each block, and second, to provide information about variations in population density within inhabited sections of each block. We use two different ancillary datasets: the National Land Cover Dataset 2011 (NLCD 2011) (Homer et al., 2015) and the National Land Use Dataset 2010 (NLUD 2010) (Theobald, 2014). Both NLCD 2011 and NLUD 2010 grids have a resolution of 30 m and are co-registered with the grid of rasterized blocks. All three grids are in the Albers Conical Equal Area (EPSG 5070) projection so each grid cell has approximately the same area.

Within the conterminous U.S. each 30 m cell in NLCD 2011 is assigned one of possible 16 land cover classes: open water (11), perennial ice/snow (12), developed, open space (21), developed, low intensity (22), developed, medium intensity (23), developed high intensity (24), barren land (31), deciduous forest (41), evergreen forest (42), mixed forest (43), shrub (52), grassland (71), pasture (81), crops (82), woody wetlands (90), herbaceous wetlands (92). The numbers in brackets are the numerical labels of land cover classes. The problem with using the NLCD as an ancillary variable

221

222

223

224

226

227

228

229

230

231



Figure 1: Decision tree showing the process of assigning a cell's ancillary class on the basis of its NLCD and NLUD classes and the population count in the block to which it belongs

for dasymetric modeling is that its classes are based on 241 surface spectral properties and thus cannot distinguish 242 between populated buildings and unpopulated buildings 243 as well as other impervious surfaces. 244

To have better information about populated vs. un-245 populated areas we also utilize the NLUD 2010. Each 246 298 cell in NLUD 2010 is assigned one of 79 land use 247 299 classes (see Table 1 in Theobald (2014)). These classes $_{300}$ 248 are divided into five broad categories: water, built-up, 301 249 production, recreation, and conservation. A built-up 302 250 category is further subdivided into residential and non-251 303 residential (commercial, industrial, institutional, trans-252 portation) subcategories. In principle, NLUD is supe-253 rior to NLCD as an ancillary variable for dasymetric 254 modeling because it relates more directly to population 255 density. Therefore, it could be argued that the dasymet-256 308 ric model should be constructed exclusively on the ba-257 sis of NLUD. However, NLUD contains artifacts due to 258 fact that it is a compilation of many different datasets of 259 varying quality and form. Thus, for the purpose of our 310 260 dasymetric model we utilize NLUD 2010 only to distin-261 311 guish between inhabited and uninhabited areas - a dis-262 312 tinction that is the most problematic while using NLCD. 313 263 Thus, we reclassify NLUD into two categories: unin- 314 264 habited (water, built-up commercial, built-up industrial, 315 265 built-up institutional (except nursing homes), built-up 316 266 transportation, mining area, and general and developed 267 317 parks) and inhabited (all other classes). 26

Information from the census blocks, NLCD 2010 and 319 269 the reclassified NLUD 2010 is combined to assign to 320 270 each grid cell one of 6 possible ancillary classes: un-321 271 inhabited (6), inhabited, vegetation (5), inhabited, high 322 272

intensity (4), inhabited, medium intensity (3), inhabited low intensity (2), and inhabited, open space (1). The number in brackets are our numerical labels for these classes.

The process of assigning these classes is illustrated by the decision tree shown in Fig. 1. Each cell is subjected to a hierarchy of predicate statements to decide its ancillary class. At the first node the cell is assigned to class 6 (uninhabited) if the block to which it belongs has population count of zero. At the second node the cell is assigned to class 6 if the cell has NLCD label 11 (open water), 12 (perennial ice/snow), or 31 (barren land). The third node tests whether the cell is assigned to inhabited or uninhabited categories according to our reclassification of NLUD. If the cell has an "inhabited" category, the fifth node assigns it to one of 5 inhabited classes (1 - open space, 2 - low intensity, 3 - medium intensity, 4 - high intensity, 5 - vegetation) based on its NLCD labels as shown in Fig. 1. If the cell has an "uninhabited" category then the fourth node assigns it to class 6 unless all cells in a given block are assigned to the "uninhabited" category which is in the conflict with the fact that the block has a nonzero population count (follow the tree back to the first node). As the census information is considered more reliable, the cell is assigned its ancillary class based on its NLCD labels as shown in Fig. 1.

Fig. 2 illustrates the construction of the ancillary data layer using six adjacent blocks in Cincinnati, OH as an example. A satellite image of the blocks is given for reference (panel A). Panel B shows the NLCD map; it can be checked alongside the image for accuracy. Panel C shows the division of the area into inhabited and uninhabited sections according to the NLUD. Finally, panel D shows the map of ancillary classes we derived from the NLCD and NLUD using the decision tree in Fig. 1.

3. Dasymetric model

The key step when constructing a dasymetric model is the establishment of the relationship between ancillary variables and population density. A significant body of literature exists on the different methods used to establish such a relationship. Their review is beyond the scope of this paper. However, for the most relevant background we refer the reader to the descriptions of models used in other projects whose aim was to construct large scale population grids: Gallego et al. (2011) discussed various models tested while computing the 100 m grid for countries in the European Union, and Stevens et al. (2015) discussed a model used to compute the 100 m grids for the WorldPop project.

318

286

287

288

289

290

291

292

293

20/

295

296

297

304

205

306



Figure 2: Construction of the ancillary layer using area consisting of six adjacent blocks in Cincinnati, OH as an example. (A) A satellite image showing the surface masked to the spatial extent of the area; six constituent blocks are indicated by orange lines. (B) Spatial distribution of NLCD classes over the area. (C) Spatial distribution of reclassified NLUD classes over the area. (D) Spatial distribution of final ancillary classes over the area.

In our model the relationship between the ancillary 342 323 variable (6 classes resulting from combining NLCD 343 324 and NLUD information) and population density is in 344 325 the form of characteristic values of population density 345 326 for each ancillary class (Mennis and Hultgren, 2006). 346 327 These values are established by sampling the population 347 328 density in blocks selected from the entire U.S. which 348 329 are (almost) homogeneous with respect to their ancil- 349 330 lary classes. For ancillary classes 1 to 4 we set the block 350 331 homogeneity threshold to 90% and for ancillary class 5 332 we set the block homogeneity threshold to 95%. 333 352

Table 1 summarizes the block samples. The sec- 353 334 ond column gives the count of homogeneous blocks 354 335 with respect to each ancillary class, the third column 355 336 gives the population count in homogeneous blocks, and 356 337 the fourth column gives the total area of homogeneous 357 338 blocks. The next three columns (fifth to seventh) give 358 339 different estimates of characteristic population density 359 340 (in people/km²) for each ancillary class. Values in the 360 341

fifth column are calculated by dividing the entire population in a given sample by the total area of all blocks in this sample. Values in the sixth column are the medians of population densities of individual blocks in a given sample. Values in the seventh column correspond to the maximum probability of a given sample probability distribution function. Only for class 4 (inhabited, high intensity) do the three estimates of population density vary considerably.

Fig. 3 shows the probability distribution functions for values of population density in each sample of homogeneous blocks. From these distributions it is clear that the density values are broadly distributed. Despite this broadness, the shapes of density distribution functions for ancillary classes 1,2,3, and 5 indicate the existence of characteristic values – values of maximum probability which correspond to the values of a sample's mean as well as its median. However, this is not the case for ancillary class 4 (inhabited, high intensity) for which the

Table 1: Characteristic population densities for six ancillary classes

Ancillary class	Blocks count	Population	Area [km ²]	Mean density	Median density	Max. prob	d [%]
uninhabited (6)	4576562	0	3001721	0	0	0	0
inhabited, open space (1)	47198	701428	902	778	850	801	4.21
inhabited, low intensity (2)	208171	5571086	2511	2219	2051	1874	11.98
inhabited, medium intensity (3)	143543	7847399	1648	4761	4060	3346	25.73
inhabited, high intensity (4)	30887	2966009	276	10743	6389	1076	58.05
inhabited vegetation (5)	753140	11795302	2153420	5	6	6	0.03



391 Figure 3: Probability distribution functions of population density val-392 ues in samples of blocks which are homogeneous with respect to ancillary class. The distribution for ancillary class 5 is shown as an inset 393 because of its vastly different scale of density values. The colors of distribution curves and the numerical labels of the curves correspond to ancillary class labels (see legend of Fig. 2).

probability peak (albeit a small one) occurs at a much 395 361 smaller density value than the sample average – there is 396 362 397 simply no good characteristic value of population den-363 398 sity associated with this ancillary class. 364

We have chosen to use the values in the fifth column 399 365 (average method) of Table 1 as sampled characteristic 366 population densities and denoted them by the symbols 367 p_i , where $i = 1 \dots 6$ are numerical labels of ancillary 368 classes. The last column in Table 1 gives the values of 369 relative population density coefficients d_i 370

$$d_i = \frac{p_i}{p_1 + \dots + p_6} \times 100, \quad i = 1, \dots 6 \tag{1}_{404}$$

The relative population density coefficients could be 371 406 interpreted as percentages of a given block population 407 372 assigned to portions of the block covered by corre- 408 373 sponding ancillary classes in a hypothetical case when 409 374 the block area is divided equally between all ancillary 410 375 376 classes. 411

According to this model the relative population den- 412 377 sity in areas belonging to the high intensity class is 378 413

about twice the relative density of areas belonging to the medium intensity class, which, in turn is about twice the relative density of the low intensity class. The relative density of the open space class is about three times lower than that of the low intensity class. Finally, the relative density of the vegetation class is about two orders of magnitude lower than that of the open space class.

A dasymetric model disaggregates population within a block to its constituent cells in proportion to a relative population density coefficient corresponding to the ancillary class associated with a cell. Denoting a given block by label x, where $x = 1, \dots, 11007989$, and denoting a given cell in this block by label *j*, we calculate the weight $W_{x,j}$ associated with this cell as:

$$W_{x,j} = \frac{d_{x,j}}{\sum_{k=1}^{6} A_{x,k} d_k}$$
(2)

where $d_{x,i}$ is the value of the relative population density coefficient in cell j of block x determined by the label of the cell's ancillary class and $A_{x,k}$ is an area (in units of cells) within cell x associated with ancillary class k. The denominator in eq. 2 is to ensure that weights over all cells in the block add to 1. The population count in cell *j* of block *x* is given as:

$$pop_{x,j} = W_{x,j} \times pop_x \tag{3}$$

where pop_x is the population in block x and $pop_{x,j}$ is the population in cell j of block x. From the properties of the weights it is clear that the sum of populations of all constituting cells is equal to the population of the block. Thus, our dasymetric model has a pycnophylactic (mass-preserving) property. Note that $pop_{x,i}$ is not an integer number and, in many areas, it is smaller than 1. Thus, referring to $pop_{x,j}$ as a population count is inappropriate. Instead, $pop_{x,j}$ should be referred to as a population density with units of people/area. Calculated values of $pop_{x,i}$ are given in units of people/900m² and our downloadable data (see section 4.3) are given in these units. For the purpose of illustration in this

394

401

402

paper, and for the online population density map, we 462 414 recalculated the values to the more customary units of 463 415 people/km² by multiplying each cell value by 9×10^{-4} 416 464 which is the area of a cell in km^2 . 417 465

Fig. 4 demonstrates the use of the dasymetric model 466 418 using the 6 blocks introduced in section 2.2 as an ex- 467 ample. Panel A shows spatial distribution of population 468 420 density inferred from block data alone; in the absence 469 421 of any additional information population density over 470 422 each block is assumed constant and equal to the number 471 423 of people in the block divided by the area of the block. 472 424 Panel B shows the spatial distribution of weight values 473 425 calculated on the basis of ancillary data (eq. 2) and panel 474 426 C shows the spatial distribution of population density 475 427 after disaggregation by the dasymetric model (eq. 3). 428 According to the model, five out of six blocks in this 429 example are characterized by a heterogeneous distribu-430 tion of population. The grid-based map (panel C) of-431 fers more specific information about where inhabitants 432 of these blocks live. The spatial precision of the grid-481 433 based map can be visually checked using a high resolu-434 482 tion image (Fig. 2A). 435

4. U.S.-wide population grid 436

Although dasymetric modeling is a well established 437 method and is relatively straightforward to apply, its 438 application to the high resolution disaggregation of 439 continental-scale areas requires an efficient computa-489 440 tional algorithm and the ability to handle big datasets. 441 The major challenge for the calculation of a 30 m dasy-442 metric model of population density for the entire con-443 terminous United States is the size of input and output 444 data. Another challenge is to provide an intuitive and 445 convenient means of reviewing and accessing the grid 446 data. 447

4.1. Computation of population grid 448

Our computational task was to disaggregate over 11 499 449 millions census blocks into over 8 billion 30 m grid 500 450 cells, and to do this in a reasonable time on a relatively 501 451 modest computer with Intel 3.4GHz, 4-core processor 502 452 and 16 GB of memory running the Linux operating sys- 503 453 tem. We use a combination of two open source software 504 packages: GRASS 7.0 (Neteler and Mitasova, 2007) 455 which is a geographical information system platform, 506 456 and R (R Development Core Team, 2008) which is a 507 457 458 programming language and software environment. In 508 addition, we also use the spgrass6 (Bivand, 2007) pack-509 459 age that allows for the efficient transport of data between 510 460 GRASS and R. 461

Computation consists of several steps that follow the methodology described in sections 2 and 3: data preprocessing (in GRASS), calculation of the population grid (in R), and post-processing (in GRASS). The input for the pre-processing step are vector data containing the boundaries of census blocks together with an attribute table as well as ancillary datasets (NLCD and NLUD) for each county separately. Pre-processing performs block rasterization and computes a single ancillary dataset from NLCD+NLUD datasets (see Fig. 1). At the end it imports the data (organized by county) to R. The pre-processing step takes 37 h. Actual calculation of the grid using the dasymetric model (see section 3) is performed in R and consists of establishing weights and disaggregating the block population according to these weights. This step takes 6 h. The post-processing step consists of exporting the grid (organized by county) to GRASS and joining grids for separate counties into one grid for the entire conterminous U.S. This step takes 18 h. Altogether, the entire computation takes 61 hours. However, the grid of weights is stored and can be reused for disaggregation of other block-level variables which are related to population, such as sex, age, and race. Thus, for example, to obtain a U.S.-wide grid of the population of African-Americans we would start from weight we have already calculated and perform only disaggregation and post-processing.

4.2. Assessment of accuracy

The accuracy of the grid can be assessed directly only if ground truth data is available. In this context the ground truth data would consist of certifiable population counts within aggregation units smaller than those used in the dasymetric model. For example, the EU population grid (Gallego, 2010; Gallego et al., 2011) was obtained by disaggregating population in communes relatively large areal aggregation units with areas much larger than 1 km². In this case, ground truth data in the form of 1 km² population grids was available for several countries in the EU. On the other hand, census blocks - the areal aggregation units we chose to disaggregate - are the smallest aggregation units available for the U.S., thus, we lack any sub-block population ground truth data to assess the accuracy of our grid directly.

Instead, we assess the accuracy of our method by calculating an additional grid based on the disaggregation of larger units - census block groups - and comparing the population of the resultant grid aggregated to blocks with the population of the blocks as given by the census. This method was used previously, including most recently by Jia et al. (2014) and Stevens et al. (2015).

483

484

485

486

487

488

490

491

492

493

495

496

497



Figure 4: Demonstration of dasymetric modeling using an area consisting of six adjacent blocks in Cincinnati, OH as an example. (A) Map of population density using only block-level data. Numbers indicate population counts for each block. (B) Spatial distribution of weight values. (C) Map of population density using grid data calculated using a dasymetric model.

535

536

537

538

539

540

541

542

543

544

Let's consider a given block group consisting of M 512 blocks. We denote the population of *m*-th block, as 513 given by block level data, by pop_m^{GT} , where the super-51 script GT indicates ground truth. We denote the popula-515 tion of *m*-th block as obtained from group-based dasy-516 metric model by pop_m^{DM} , where the superscript DM in-517 dicates dasymetric model. Jia et al. (2014) following 518 Eicher and Brewer (2001) used two quantities to assess 519 the accuracy of their methods: 520

$$\text{RMSE} = \sqrt{\frac{1}{M} \sum_{m=1}^{M} (pop_m^{GT} - pop_m^{DM})^2} \qquad (4)$$

521 and

$$CV = \frac{RMSE}{\frac{1}{M} \sum_{m=1}^{M} pop_m^{GT}}$$
(5)

where RMSE is the root square mean error and CV is 545 522 the coefficient of variance. RMSE expresses the devi-523 546 ation (in the number of people) of the model from the 547 524 ground truth in absolute terms, whereas CV expresses 548 525 this deviation relative to the population. Both quantities 549 526 pertain to a single block group. Calculating the mean (or 550 527 median) of these quantities over all block groups in the 551 52 region covered by the grid assesses the overall accuracy 552 529 of the grid. 530 553

We first calculate the statistics of RSME and CV for our grid in order to compare them to the analogous statistics calculated by Jia et al. (2014) for their 30 m resolution population grids covering Alachua county in 557

Table 2:	Accuracy	assessment	using	RMSE and CV

Tuble 2. Theedracy assessment using fulled and C v							
Grid	Mean RMSE	Mean CV	Median CV				
Conterminous U.S.	43.17	0.97	0.78				
Alachua	63.12	1.29	1.21				
Alachua NLCD	73.26	1.36	1.30				
Alachua parcels	63.96	1.20	0.88				

the state of Florida. They calculated two grids (disaggregated from 2010 census blocks) using two different ancillary datasets: NLCD, and 2010 tax parcel data. Parcel data is considered better ancillary data for dasymetric modeling than land cover because it relates more directly to population count (Jia et al., 2014). However, it is only available (in various degrees of completeness) for nineteen states in the U.S. (Stage and VonMeyer, 2006) and thus cannot constitute a base for a U.S.–wide population grid.

The first row in Table 2 shows the mean value of RMSE, the mean value of CV, and the median of CV, respectively. These statistics were calculated over all block groups in the conterminous U.S. using our dasymetric model. For the remaining three rows in Table 2 statistics were calculated only over block groups in Alachua county. The second row shows values calculated using our model, the third row shows values calculated using Jia et al. (2014) model utilizing NLCD as an ancillary variable, and the fourth row shows values calculated using the Jia et al. (2014) model utilizing tax parcel data as an ancillary variable. Statistics for Alachua county indicate that our grid (restricted to

Alachua county) has higher accuracy than the Jia et al. 558 (2014) model based on NLCD despite using values of 559 characteristic population densities derived from nation-560 wide block samples. This can be attributed to our more 561 advanced dasymetric model which utilizes NLCD as 562 well as NLUD, and to having six ancillary classes in-563 stead of three. In comparison to our model the Jia et al. 56 (2014) model based on parcel data has a slightly higher 565 value of $\langle RMSE \rangle$, a slightly lower value of $\langle CV \rangle$, and a 566 lower value of median CV value. Thus, the population 567 grid (within the Alachua county) constructed using tax 568 parcels as ancillary data has somewhat higher accuracy 569 than our grid constructed using a combination of NLCD 570 and NLUD but this difference is small. Assessment of 57 our grid over the entire conterminous U.S. yields supe-572 rior accuracy in comparison with the part restricted to 573 Alachua county and in comparison with both grids con-574 structed by Jia et al. (2014). 575

Although the accuracy measures discussed above are 576 useful for the comparison of different dasymetric mod-577 els, they do not indicate to a user the degree of accu-578 racy that can be expected from a grid. If a user uti-579 lizes the grid to estimate the population count of a sub-580 block areal unit, what is an uncertainty of this esti-581 mation? To provide this information we assume that 582 uncertainty indicators calculated for blocks using the 583 dasymetric model obtained by disaggregation of block 584 groups are valid for sub-block units when using a dasy-585 metric model which disaggregate blocks. 586

We calculate the statistics of relative errors over all 611 587 612 inhabited blocks in the conterminous U.S. Note that due 588 613 to the way we construct our ancillary classes (see Fig. 1) 589 614 uninhabited blocks have zero population error in our 590 615 model and could be excluded from statistics. The rel-591 616 ative error δ_m of the population count for block *m* is the 592 617 absolute error divided by the magnitude of the ground 593 truth value. 594

$$\delta_m = \frac{|pop_m^{GT} - pop_m^{DM}|}{pop_m^{GT}} \tag{6}$$

The value of δ_m has a simple interpretation – if multi-595 plied by 100 it expresses the overestimation or underes- 624 596 timation of the number of people as a percentage of the 625 597 actual population of the block. As the distribution of 626 598 the values of δ_m over all inhabited blocks is skewed, the 627 599 usual statistical indicators (mean and standard variation) 628 600 are not providing useful information, instead a median 629 601 602 (equal to 0.44) provides a robust estimation of the "ex- 630 pected" value of δ_M and the median absolute deviation 631 603 (equal to 0.4) provides a robust estimation of spread in 632 604 the values of δ_m around the expected value. Thus, when 633 605



Figure 5: Two-dimensional histogram of all inhabited blocks with respect to their population density and the value of relative error between its modeled and actual populations. Number of blocks in each bin of the histogram is color-coded according to the legend. Note logarithmic scales for all variables.

using our grid to estimate a population in a sub-block area, a user can expect, on average, 44% uncertainty in the population count.

We can provide further information on the uncertainty of grid-based population estimations by calculating a two-dimensional histogram (shown in Fig. 5) of blocks with respect to their population density and relative error. Note that both population density (horizontal axis) and relative error (vertical axis) are shown in logarithmic scale due to the orders of magnitude variations in their values. The histogram has a final number of bins with each bin represented by a small square in Fig.5. The color of a bin carries information on the number of blocks having the population density and relative error as indicated by the bin's coordinates. Note that the block count legend is also logarithmic, 83% of all blocks are in the three top block counts categories indicated by pink, darker red, and lighter red colors.

One way to use Fig. 5 to obtain information about population count uncertainty is to first select a value of population density for an area of interest. Values of block counts along the vertical line corresponding to the selected value of density combine into distribution of error values for these blocks. The peak of this distribution is located at the expected value of an error for blocks in this population density range and a deduced shape of the distribution informs about spread of error values around the expected value. Thus, assuming that histogram ob-

607

609

610

618



Figure 6: Screenshot of SocScape after population density data layer 676 has been selected, the map zoomed to show the area around Chicago, 677 IL. The download tool is selected from the navigation panel, and download region is indicated by a user.

tained for blocks is also valid for sub-block areas a user 634 interested in an area having population density of about 635 3000 people/km² finds that the error is between 0.23 and 636 0.45. On the other hand, a user interested in areas hav-637 ing population density of about 100 people/km² finds 638 that the error is most likely to be about 0.67. 639

4.3. Data access 640

We provide convenient access to the 2010 population 691 641 grid via SocScape (Social Landscape) – a GeoWeb ap-642 plication designed for exploration and data distribution 693 643 of population density and racial diversity grids. Soc- 694 644 Scape is available at http://sil.uc.edu. Upon launching 695 645 SocScape shows a background map of the United States 696 646 and a menu to select data. When "Population density 697 647 2010" is selected from the menu a map of population 698 density appears categorized to eleven bins represented 699 649 by different colors (the map legend is accessible from 700 650 the navigation panel). It is important to differentiate be-651 tween the map of population density, which is intended 702 652 only for online exploration, and actual (not categorized) 703 653 data that can be downloaded once a user decides on an 704 654 area of interest. 655

Data in the population grid corresponds to population 706 656 density and has units of people per cell. It can be down-707 657 loaded by an area of interest. To download the data a 658 user has to zoom into a general area of interest and select 709 659 660 the download tool from the navigation panel. When the 710 download tool appears the "Population density 2010" 661 data layer needs to be selected. Next, the user indicates 712 662 a specific region of interest by dragging the mouse to 663

draw a rectangle. Pressing on the GeoTIFF button in the download tool will start the download. The data is provided in geotiff format.

5. Grid-based versus blocks-based population maps

What differences can one expected when mapping population using population density (an intensive variable defined on a regular grid) versus a population aggregated to units (an extensive variable defined on census blocks). To point out differences between the two approaches we selected two locations as examples, first - a highly urbanized area of San Francisco CA, and second - a rural area centered on the Lake Loramie State Park, OH.

Fig. 7A shows a map of population density in San Francisco calculated from census blocks; density is uniform over each block as it was calculated by dividing the population in a block by its area. Block boundaries are not shown because the size of the blocks in this area are very small and the lines representing their boundaries would obscure the map at the figure's level of resolution. Fig. 7B shows a map of population density as represented by our grid. Overall, the two maps are similar but there some differences that can be summarized into two categories: (a) the block-based map is more consolidated and thus appears to have less detail, and (b) the grid-based map shows uninhabited areas which appear as inhabited on the block-based map.

The first difference stems directly from the fact that the dasymetric model disaggregates blocks into subblock cells whose densities may vary. In general this results in superior spatial resolution for the grid. However, we need to keep in mind the uncertainty of the dasymetric model. Consider areas immediately south and north of Golden Gate Park (a prominent rectangle elongated along the west-east axis). Inspection of these areas using a high resolution image or map (for example those available in Google Maps) reveals that they are characterized by a grid layout of streets with houses filling the grid. This is not captured by the blockbased map because the blocks contain both houses and streets. With 30 m resolution one would expect that the grid-based map would show a white-purple pattern corresponding to streets (uninhabited) and houses (inhabited with high density). Instead, we observe a pattern consisting of purple (high population density) and lightpurple (smaller population density) colors which is due to the fact that the NLCD in these areas does not recognize narrow streets and interprets the Landsat image as a mosaic of high and medium intensity developed land cover classes; the NLUD also do not delineate streets in

711

678

679

680

681

682

685

686

687

688

689

690

701



Figure 7: Comparison of population maps for San Francisco (CA). (A) The block-based map. (B) The grid-based map. Boundaries of census blocks are not show for clarity.

742

743

745

746

747

748

749

750

751

752

753

766

756

these areas. Thus, although the grid correctly indicates 738 714 the heterogeneity of population density in these areas, 739 715 the spatial scale of streets is too small for them to be rec-740 716 ognized as uninhabited areas. Still, there is some gain 741 717 in information with respect to the block-based map. 718

The second difference is a direct result of the mod-719 ifiable areal unit problem, entire blocks in the block-720 based map are marked by colors corresponding to low 721 or medium values of population density even if they are 722 mostly uninhabited but have small portions occupied by 723 housing. Inspection of a high-resolution image or map 724 reveals that most areas on the block-based map shown 725 in yellow, light brown, or red colors are really predomi-726 nantly uninhabited. 727

Fig. 8A shows a map of population density in a ru-728 ral site centered on the Lake Loramie State Park, OH 754 729 calculated from census blocks; density is uniform over 730 each block as it was calculated by dividing the popula-731 tion in a block by its area. A high resolution image of 732 this site (Google Earth) reveals the presence of Lake Lo-733 734 ramie (center) and three small towns (the western part 757 of the site), but most of the site is an agricultural land-735 scape crossed by a grid-like network of secondary roads. 758 736 Farmhouses are located predominantly at intersections 759 737

of the roads, leaving most of the land uninhabited.

The block-based map does not reflect the real distribution of population because it assigns homogeneous density to predominantly uninhabited blocks. As a result most of the map is shown in a brown color corresponding to a population density of 10-50 people/km². This is a relatively small value of density but still the reality is that these areas are uninhabited except at the locations of individual farmhouses. Fig. 8B shows the map of population density as represented by our grid. This map does not completely eliminate the inaccurate impression about the character of population distribution at this site but it significantly alleviates the problem by concentrating population in farmhouses and along the roads leaving the rest of the countryside either uninhabited or with a negligibly small density of population. Fig. 8B must be magnified in order to see small red dots indicating the population concentration along the secondary roads and at individual locations.

6. Discussion and conclusions

The purpose of the work presented in this paper, was to deliver the best possible nationwide population grid



Figure 8: Comparison of population maps for Lake Loramie State Park (OH) site. (A) The block-based map; block boundaries are shown as black lines. (B) The grid-based map.

801

that can be constructed using readily available public 794 760 domain data. The resultant product, which we call 795 761 SocScape–30, is a 30 m resolution grid carrying values 796 762 of residential (nighttime) population densities in 2010. 763 The grid is freely distributed through the web using the 797 764 SocScape GeoWeb application. The availability of this 798 765 resource should make population data more accessible 799 766 and thus more utilized. 800 767

We decided to combine land cover (NLCD) and land 802 768 use (NLUD) datasets into a single ancillary variable to 803 769 guide dasymetric modeling - the technique on which 804 770 our grid is based. It may be asked why we did not uti-805 771 lize more ancillary datasets such as tax parcel data, road 772 network data, the density of points of interest, topogra-807 773 phy, light emission etc. There are three reasons behind 808 774 our choice. First, unlike other parts of the world, in 809 775 the U.S. available land cover and land use datasets are 810 776 highly reliable and the smallest census units (blocks) 811 777 are small making additional ancillary information re- 812 778 dundant or unnecessary. Second, apart from the road 779 network (which we claim is redundant when NLUD is 814 780 utilized) land cover and land use grids are the only read- 815 781 ily available ancillary data that are consistent over the 816 782 entire U.S. Tax parcel data – potentially useful ancillary 817 783 information – is only available for some states (see sec- 818 784 tion 4.2) and each state or even county releases this data 819 785 in its own format making the consistency of tax parcel 820 786 data an issue. Our uncertainty estimates (section 4.2), 821 787 admittingly performed only for a single county, indicate 822 788 that using ancillary data based on the combination of 789 823 790 NLCD and NLUD yields a grid that is only slightly less 824 accurate than a grid based on disaggregation using tax 825 791 parcel data. Point of interest data is equally inconsistent 826 792 on the continental scale. Third, land cover and land use 827 793

datasets are available in a convenient grid format making data pre-processing more efficient and free from potential artifacts.

We use a relatively simple dasymetric model instead of seemingly more advanced models based on supervised machine learning (Stevens et al., 2015) because our model uses only a single ancillary variable. The model is based on nation-wide statistics rather then on a series of local statistics which could potentially capture a non-stationarity of the relationship between population density and an ancillary variable (Lo, 2008; Gallego, 2010; Schroeder and VanRiper, 2013). We selected this model from among the three we have calculated. The other two models attempted to address the non-stationarity issue. For the second model we divided the U.S. into five zones following the United States Department of Agriculture (USDA) Rural-Urban Continuum Codes for U.S. counties: rural areas, small town, micropolitan areas, metropolitan statistical areas (MSAs), and MSAs with population > 1,000,000 people. For each zone separately we calculated a dasymetric model as described in section 3. As expected, the values of relative population density coefficients vary somewhat from one zone to another but the accuracy of such model, as measured by the mean value of CV calculated over the entire conterminous U.S., is not higher than that of our default model. In addition, the zoned model contains artifacts as the population density is discontinuous at the boundaries between the zones. In our third model we fitted characteristic values of population density separately for each county or census tract (like in Gallego (2010) who deployed such a technique for Nomenclature of Territorial Units for Statistics (NUTS) - very large territorial regions) but we found that tracts

or counties are too small to have a statistically valid 880 828 sample of homogeneous blocks. Thus, overall, we 881 829 deemed the default model to be the best choice for our 882 830 purpose. 831 883

We have performed a comprehensive assessment of 884 832 the accuracy of the method used to obtain our grid (sec- 885 833 tion 4.2). We estimate an overall relative error to be 886 834 44%, which may appear to be large. However, it is at the 835 lower limit of errors estimated for other methods used 888 836 to create state-of-the-art population grids. The error of 889 837 the EU 100 m grid (Gallego et al., 2011) is estimated to 890 838 be between 40% to 105% depending on the dasymetric 891 839 model used and the country for which the assessment 892 840 was done. The error for the WorldPop population grids 893 is estimated to be between 39% and 91% for the newest 842 models (Stevens et al., 2015) or between 46% to 120% 843 895 for older models (Gaughan et al., 2013) depending on 844 896 the country. It is important to note that each project uses 845 a slightly different methodology to assess accuracy, but 846 conclusions are similar - dasymetric model has an ex-897 847 pected uncertainty of about 40-100%. We also provide 848 the means for a more specific estimation of expected 899 error (Fig. 5) based on additional knowledge about the 850 population density of the area of interest. The largest 851 902 source of error stems from the blurred relationship be-852 tween land cover classes and population density. Using 904 853 land use classes would help with this problem but the 854 quality of U.S.-wide land use data is not sufficient to be 855 907 utilized in a role other than for delineation of uninhab-856 908 ited areas. It is important to realize that the expected 909 857 910 error estimates give the difference between a predicted 858 and an actual number of people in a area of interest. 859 912 It does not comment on the spatial precision of delin-860 913 eation between inhabited and uninhabited areas. In our 914 861 915 grid this delineation follows the land use data and is ex-862 916 pected to be fairly accurate. 863 917

As we mentioned in section 4.1 weights (eq. 2) cal-918 culated on the basis of our model can be used to disag-865 gregate other census block-level (SF1) variables which 866 represent population segments; examples include sex, 867 age and race. Such a disaggregation will not be able 923 868 to account for possible sub-block heterogeneities in the 869 proportion of each population segment to the total popu-870 lation, this will remain fixed throughout each block, but 927 871 it will give a more accurate spatial location of each pop-928 872 ulation segment by keeping it away from uninhabited 873 areas and making an adjustment in line with the sub-874 block overall population density. One immediate appli-875 876 cation is the construction of a 2010 version of diversity maps like those introduced in Dmowska and Stepinski 877 (2014) and analyzed in Dmowska and Stepinski (2016) 878 for years 1990 and 2000. 879

Future plans call for recalculation of the 1990 and 2000 editions of U.S. population grids from 90 m resolution obtained by disaggregating SEDAC 1 km grid (Dmowska and Stepinski, 2014) to 30 m resolution using the technique presented in this paper. The major promise of having grids for various years is the ability to assess spatio-temporal population change. For such grids not to lead to discovery of spurious change they must be constructed using compatible ancillary datasets. Our plans call for us to make available through Soc-Scape not only the best possible population grid for 2010 as described in this paper - but also for compatible, but not necessarily the best possible - population grids for 1990, 2000, and 2010 for the purpose of analyzing spatio-temporal change.

Acknowledgments. This work was supported by the University of Cincinnati Space Exploration Institute.

References

- Balcik, B., Beamon, B. M., 2008. Facility location in humanitarian relief. International Journal of Logistics 11 (2), 101-121.
- Benn, H. P., 1995. Synthesis of transit practice 10: bus route evaluation standards. Tech. rep., Transit Cooperative Research Program, Transportation Research Board, National Research Council, Washington, DC
- Bhaduri, B., Bright, E., Coleman, P., Urban, M. L., 2007. Land-Scan USA: a high-resolution geospatial and temporal modeling approach for population distribution and dynamics. GeoJournal 69 (1-2), 103-117.
- Bivand, R., 2007. Using the R GRASS Interface: Current Status. OSGeo Journal 1, 36-38.
- Briggs, D. J., Gulliver, J., Fecht, D., Vienneau, D. M., 2007. Dasymetric modelling of small-area population distribution using land cover and light emissions data. Remote sensing of Environment 108 (4), 451-466.
- Chen, K., McAneney, J., Blong, R., Leigh, R., Hunter, L., Magill, C., 2004. Defining area at risk and its effect in catastrophe loss estimation: A dasymetric mapping approach. Applied Geography 24, 97-117.
- Deng, C., Wu, C., Wang, L., 2010. Improving the housing-unit method for small-area population estimation using remote-sensing and GIS information. International Journal of Remote Sensing 31(21), 5673-5688.
- Dmowska, A., Stepinski, T. F., 2014. High resolution dasymetric model of U.S demographics with application to spatial distribution of racial diversity. Applied Geography 53, 417-426.
- Dmowska, A., Stepinski, T. F., 2016. Mapping changes in spatial patterns of racial diversity across the entire United States with application to a 1990-2000 period. Applied Geography 68, 1-8.
- Dobson, J. E., Brlght, E. A., Coleman, P. R., Worley, B. A., 2000. LandScan: a global population database for estimating populations at risk 66 (7), 849-857.
- Eicher, C. L., Brewer, C. A., 2001. Dasymetric mapping and areal interpolation: Implementation and evaluation. Cartography and Geographic Information Science . 28, 125-138.
- Flowerdew, R., Green, M., 1992. Developments in areal interpolation methods and GIS. Annals of Regional Science 26, 67-78.
- Gallego, F., Batista, F., Rocha, C., Mubareka, S., 2011. Disaggregating population density of the European Union with CORINE land

919

920

921

925

931

932

933

934

- cover. International Journal of Geographical Information Science 1003
 25 (February 2015), 2051–2069. 1004
- Gallego, F. J., 2010. A population density grid of the European Union.
 Population and Environment 31 (6), 460–473.
- Gaughan, A. E., Stevens, F. R., Linard, C., Jia, P., Tatem, A. J., 2013. 1007
 High resolution population distribution maps for Southeast Asia in 1008
 2010 and 2015. PLoS One 8(2), e55882. 1009
- 945Gleick, P. H., 1996. Basic water requirements for human activities: 1010946Meeting basic needs. Water international 21 (2), 83–92.1011
- Goodchild, M., Anselin, L., Deichmann, U., 1993. A framework for 1012
 the areal interpolation of socioeconomic data. Environment and 1013
 Planning, A 25, 383–397. 1014
- Goodchild, M., Lam, N., 1980. Areal interpolation: A variant of the 1015
 traditional spatial problem. Geo-Processing 1, 297–312.
- Hay, S. I., Noor, A. M., Nelson, A., Tatem, A. J., 2005. The accuracy 1017
 of human population maps for public health application. Tropical 1018
 Medicine & International Health 10(10), 1073–1086.

Holt, J. B., Lo, C. P., Hodler, T. W., 2004. Dasymetric estimation 1020

- 956of population density and areal interpolation of census data. Car-1021957tography and Geographic Information Science 31(2), pp.103-121.102395831(2), 103-121.1023
- Homer, C. G., Dewitz, J. A., Yang, L., Jin, S., Danielson, P., Xian, G., 1024
 Coulston, J., Herold, N. D., Wickham, J. D., Megown, K., 2015. 1025
 Completion of the 2011 National Land Cover Database for the 1026
 conterminous United States-Representing a decade of land cover 1027
 change information. Photogrammetric Engineering and Remote 1028
 Sensing 81(5), 345–354. 1029
- Jia, P., Gaughan, A. E., 2016. Dasymetric modeling: A hybrid ap proach using land cover and tax parcel data for mapping population
 in Alachua County, Florida. Applied Geography 66, 100–108.
- Jia, P., Qiu, Y., Gaughan, A. E., 2014. A fine-scale spatial popula- 1033
 tion distribution on the High-resolution Gridded Population Sur- 1034
 face and application in Alachua County, Florida. Applied Geogra- 1035
 phy 50, 99–107. 1036
- Kar, B., Hodgson, M. E., 2012. A Process Oriented Areal Interpolation Technique: A Coastal County Example. Cartography and 1038 Geographic Information Science 39 (1), 3–16.
- Langford, M., Unwin, D. J., 1994. Generating and mapping popu- 1040
 lation density surfaces within a geographical information system. 1041
 The Cartographic journal 31 (1), 21–6. 1042
- Linard, C., Gilbert, M., Snow, R. V., Noor, A. M., Tatem, A. J., 1043
 2012. Population Distribution, Settlement Patterns and Accessibil ity across Africa in 2010. PLoS ONE 7, e31743.
- Linard, C., Gilbert, M., Tatem, A. J., 2011. Assessing the use of global 1046
 land cover data for guiding large area population distribution mod elling. GeoJournal 76, 525–538.
- Lloyd, C. D., 2014. The Modifiable Areal Unit Problem, in Exploring 1049
 Spatial Scale in Geography. John Wiley & Sons, Ltd, Chichester, 1050
 UK. 1051
- Lo, C. P., 2008. Population estimation using geographically weighted 1052
 regression. GIScience & Remote Sensing 45(2), 131–148.
- Lu, Z., Im, J., Quackenbush, L., Halligan, K., 2010. Population estimation based on multi-sensor data fusion. International Journal of 1055 Remote Sensing 31 (21), 5587–5604.
- Lung, T., Lübker, T., Ngochoch, J. K., Schaab, G., 2013. Human pop- 1057
 ulation distribution modelling at regional level using very high res- 1058
 olution satellite imagery. Applied Geography 41, 36–45. 1059
- Maantay, J., Maroko, A., 2009. Mapping urban risk: Flood hazards, 1060
 race, & environmental justice in New York. Applied Geography 1061
 29, no. 1 (2009): 29 (1), 111–124.
- ³⁰⁷ ³⁰⁸ ³⁰⁹ ³⁰⁹
- lation distribution in the urban environment: The Cadastral-based 1064
 Expert Dasymetric System (CEDS). Cartography and Geographic 1065
- 1001 Information Science 34 (2), 77–102. 106
- 1002 Martin, D., Williams, H. C., 1992. Market-area analysis and accessi- 1067

bility to primary health-care centres. Environment and Planning A 24 (7), 1009–1019.

- Mennis, J., 2003. Generating surface models of population using dasymetric mapping. The Professional Geographer 55 (1), 31–42.
- Mennis, J., 2009. Dasymetric mapping for estimating population in small areas. Geography Compass 3 (2), 727–745.
- Mennis, J., Hultgren, T., 2006. Intelligent dasymetric mapping and its application to areal interpolation. Cartography and Geographic Information Science 33 (3), 179–194.
- Mitsova, D., Esnard, A. M., Li, Y., 2012. Using enhanced dasymetric mapping techniques to improve the spatial accuracy of sea level rise vulnerability assessments. Journal of Coastal Conservation 16 (3), 355–372.
- Murray, A. T., Davis, R., Stimson, R. J., Ferreira., L., 1998. Public transportation access. Transportation Research Part D: Transport and Environment 3(5), 319–328.
- Neteler, M., Mitasova, H., 2007. Open source GIS: a GRASS GIS approach, 3rd Edition. Springer, New York.
- Pattnaik, S. B., Mohan, S., Tom, V. M., 1998. Urban bus transit route network design using genetic algorithm. Journal of transportation engineering 124, no. 4 (1998): 124(4), 368–375.
- Petrov, A., 2012. One hundred years of dasymetric mapping: back to the origin. Cartographic Journal 49 (3), 256–264.
- R Development Core Team, 2008. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria.
- Reibel, M., Bufalino, M. E., 2005. Street-weighted interpolation techniques for demographic count estimation in incompatible zone systems. Environment and Planning A 37, 127–139.
- Ruther, M., Leyk, S., Buttenfield, B. P., 2015. Comparing the effects of an NLCD-derived dasymetric refinement on estimation accuracies for multiple areal interpolation methods. GIScience & Remote Sensing,52(2), pp.158-178. 52(2), 158–178.
- Schroeder, J. P., 2007. Target-Density Weighting Interpolation and Uncertainty Evaluation for Temporal Analysis of Census Data. Geographical Analysis 39(3), 311–335.
- Schroeder, J. P., VanRiper, D. C., 2013. Because Muncie's Densities Are Not Manhattan's: Using Geographical Weighting in the ExpectationMaximization Algorithm for Areal Interpolation. Geographical analysis 45(3), pp.216-237. 45(3), 216–237.
- Smith, S. K., Nogle, J., Cody, S., 2002. A regression approach to estimating the average number of persons per household. Demography 39 (4), 697–712.
- Sorichetta, A., Hornby, G. M., Stevens, F. R., Gaughan, A. E., Linard, C., Tatem, A. J., 2015. High-resolution gridded population datasets for Latin America and the Caribbean in 2010, 2015, and 2020. Sci. Data 2, 150045.
- Sridharan, H., Qiu, F., 2013. A Spatially Disaggregated Areal Interpolation Model Using Light Detection and Ranging-Derived Building Volumes. Geographical Analysis 45, no. 3 (2013): 45 (3), 238– 258.
- Stage, D., VonMeyer, N., 2006. An assessment of parcel data in the United States – 2005 survey results. Tech. rep., In Federal Geographic Data Subcommittee on Cadastral Data.
- Stevens, F. R., Gaughan, A. E., Linard, C., Tatem, A. J., 2015. Disaggregating census data for population mapping using Random forests with remotely-sensed and ancillary data. PloS One 10(2), e0107042.
- Su, M. D., Lin, M. C., Hsieh, H. I., Tsai, B. W., Lin, C. H., 2010. Multi-layer multi-class dasymetric mapping to estimate population distribution. Science of the total environment 408 (20), 4807–4816.
- Tenerelli, P., Gallego, J. F., Ehrlich, D., 2015. Population density modelling in support of disaster risk assessment. International Journal of Disaster Risk Reduction 13 (2015): 13, 334–341.
- Theobald, D. M., 2014. Development and applications of a compre-

- hensive land use classification and map for the US. PloS one 9 (4),e94628.
- Thieken, A. H., Müller, M., Kleist, L., Seifert, I., Borst, D., Werner,
 U., 2006. Regionalisation of asset values for risk analyses. Natural
- Hazards and Earth System Science 6 (2), 167–178.
 Ural, S., Hussain, E., Shan, J., 2011. Building population mapping
- with aerial imagery and GIS data. International Journal of Applied
 Earth Observation and Geoinformation 13 (6), 841–852.
- Vinkx, K., Visee, T., 2008. Usefulness of population files for estimation of noise hindrance effects. In: ICAO Committee on Aviation
 Environmental Protection. CAEP/8 Modelling and Database Task
- 1079 Force (MODTF). 4th Meeting. Sunnyvale, USA. pp. 20–22.
- Voss, P. R., Long, D. D., Hammer, R. B., 1999. When census geog raphy doesn't work: Using ancillary information to improve the
 spatial interpolation of demographic data. Tech. rep., Center for
 Demography and Ecology, University of Madison-Wisconsin.
- Weber, N., Christophersen, T., 2002. The influence of nongovernmental organisations on the creation of Natura 2000 during the European Policy process. Forest policy and economics 4(1), 1-12.
- Wright, J. K., 1936. A method of mapping densities of populatrion
 with Cape Cod as an example. Geographical Review 26 (1), 103–
 110.
- Zandbergen, P. A., 2011. Dasymetric mapping using high resolution
 address point datasets. Transactions in GIS 15 (s1), 5–27.