
Unsupervised Detection of Metropolitan Areas

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 City boundaries can be crisp or fuzzy depending on the effort that local governments put
2 into digital mapping. The concept of a metropolitan area is even fuzzier than the concept
3 of a city. This paper presents an unsupervised algorithm to detect metropolitan areas from
4 geographical data that is dense in urban areas and sparse in rural areas. As an example, we
5 detect metropolitan areas for the UK using the density of crime. It is possible to tessellate
6 Earth using metropolitan areas by associating to each point the metropolitan area closest to
7 it. This can potentially aid geospatial algorithms by supplying a smart divide and conquer
8 strategy.

9 1 Introduction

10 Geospatial algorithms are often designed with the human in mind. A few examples include (a) Resource
11 allocation problems, such as determining the optimal locations for police stations, (b) Recommendation
12 problems, such as ranking the most popular tourist destinations, and (c) Estimation algorithms, such as
13 predicting real estate prices. Ideally, these problems would be solved per city, however, humans interact with
14 their environment in ways that cut across city boundaries. In one extreme, a sprawling metropolitan area has
15 no clear boundary where the “urban” ends and the “rural” starts. In another extreme, many cities can be close
16 enough together so as to form a continuum of urban landscape that can stretch hundreds of miles. When one
17 wishes to solve a problem that is local to a metropolitan area, it is not clear how to define the boundary of that
18 metropolitan area.

19 It is natural to wish to cut between metropolitan areas along lines which will least interfere with the problem
20 we are trying to solve. For the police station allocation example, this might be interpreted as a cut between
21 cities along a contour line that has minimal crime. For the recommendation example, this might be interpreted
22 as a cut along a contour line that is farthest from tourist attractions or from hotels. For the prediction
23 example, this might be interpreted as a cut along a contour line where building density is minimal.

24 We present in this paper an algorithm to find metropolitan areas. The input is any geospatial variable that can
25 be represented as a density per latitude and longitude. For the police station case, this can be the density of
26 crime; for the recommendation case, this could be the density of tourist attractions or hotels; for the estimation
27 case, this could be the density of buildings. The algorithm cuts between metropolitan areas along contour lines
28 where the variable is least dense.

29 This algorithm is based on a technique from the field of image processing known as blob detection [ref].
30 Specifically, Lindeberg’s blob detection algorithm [ref]. In this context, the density of, say, crime forms blobs
31 when visualized as a heat map. The blobs naturally correspond to metropolitan areas. We use smoothing to
32 introduce a minimal size for a blob to be considered a metropolitan area.

33 2 Metropolita area detection

34 In this section we describe our metropolitan area detection algorithm. The input is a list of latitudes and
35 longitudes. We will explore in this section a specific example: a list of locations of crimes in the UK [ref].

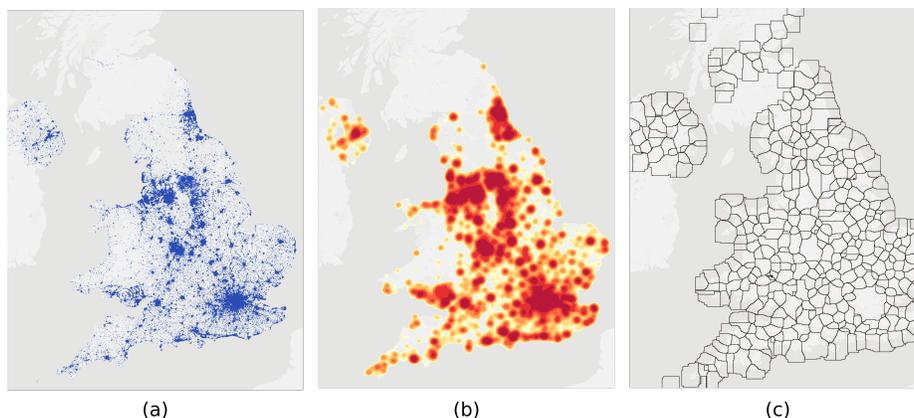


Figure 1: Metropolitan area detection applied to crimes in the UK. (a) Locations of 2 million crime events reported in January–March 2018. (b) Smoothed heat map of crime, using Gaussian smoothing of radius $\sigma = 4$ km. (c) Metropolitan areas detected by Lindeberg’s blob-detection algorithm. We plot lines at the boundary between pixels which do not share the same metropolitan area index.

36 We start by projecting the points into a local Cartesian coordinate system using the sinusoidal projection [ref]
 37 centered around a point $(\phi_0, \lambda_0) = (53^\circ N, 2^\circ W)$ which is roughly in middle of the UK. We do this so that we
 38 can use image processing techniques, such as Gaussian smoothing [ref], in the Euclidean plane. It is possible
 39 to do Gaussian smoothing on a sphere [ref] but when dealing with problems that are smaller than a large
 40 continent, a local Cartesian approximation works fine.

41 After projecting the points to a local Cartesian coordinate system, we pixelate the UK using a grid spacing
 42 of 1 km. This is measured at (ϕ_0, λ_0) and varies with latitude on a sphere, but is fixed in the Cartesian
 43 approximation. Within each pixel we count the occurrences of crime. We smooth this image using a Gaussian
 44 filter of radius σ . If two blobs are separated by a distance much smaller than this, they will be smoothed into
 45 the same blob.

46 Lindeberg’s blob detection algorithm is run on this image to enumerate the blobs and mark the contour lines
 47 that delineate a blob from its neighbors. The process is to iterate on pixels in decreasing order of crime density.
 48 If a pixel has no higher-density neighbors, it is assigned a new “metropolitan area index.” If the pixel has a
 49 neighbor with a higher density of crime, it is assigned the metropolitan area index of the neighboring pixel
 50 with highest density. Some care must be taken when there is a plateau of equal density: first expand to fill the
 51 plateau, then the entire plateau obtains a new index or copies the index of its highest-density neighboring pixel.

52 An example is shown in Figure [ref]. We downloaded UK crime data for the months of January–March 2018
 53 from [ref], dropping crimes which did not have a latitude or longitude reported. There were 2 million crimes
 54 with valid location reported in this date range. As points, they are shown in [ref]. As a smoothed heat map
 55 they are shown in [ref]. Finally, the output of the blob detection step is shown in [ref].

56 3 Comparison across datasets and smoothing parameters

57 We wish to get a feeling for the stability of this algorithm, and the robustness of the term “metropolitan area”
 58 as derived from the algorithm discussed in this paper. To that end we compare the output of the algorithm
 59 across different datasets, and with different smoothing parameters. We will compare across three open datasets:
 60 (1) The UK crime dataset described in the previous section, (2) An index of place names in the UK [ref],
 61 maintained by the office for national statistics, and (3) A list of all highway nodes in Open Street Maps [ref]
 62 for the UK [ref].

63 The term “metropolitan area“ also depends on the smoothing parameter σ . Increasing σ will decrease the
 64 number of metropolitan areas detected, and will increase the average area of each metropolitan area. We will
 65 compare the three values σ : 5km, 10km, and 20km.

66 There is a third parameter discussed in Section [ref], the grid size. Here, smaller is better, with the only
 67 tradeoff being that smaller grid sizes increase the computational time needed for the algorithm to run. For the
 68 visualizations created in this paper, a grid size of 1 km provided sufficient detail, and the algorithm ran in a



Figure 2: Comparison of our metropolitan area detection algorithm across different datasets and smoothing parameters. We show the boundaries between metropolitan areas for each of three datasets discussed in the text, and for each of $\sigma = 5, 10, 20$ km. We also overlay the figures across datasets for each value of σ to show the stability of metropolitan areas with respect to the dataset.

69 few minutes for each of the three datasets on an Intel i7 laptop. The algorithm was implemented in python and
 70 we made no attempt to optimize it.

71 The result of the comparison is shown in Figure [ref]. The three datasets are labeled “crime”, “places”, and
 72 “roads”, respectively. We show the boundaries between metropolitan areas for each dataset and for each value
 73 of σ . We also overlay the plots across datasets for each value of σ to show that the metropolitan areas are
 74 stable across datasets. The crime dataset had 2 million rows, the places dataset had 87 thousand rows, and the
 75 roads dataset had 28 million rows.