

Postcode Anonymity Analysis

Adam Hardy

May 2021

1 Introduction

In the US, 5-digit Zip codes are usually rounded to 3-digits when anonymising healthdata, so knowledge of the Zip code doesn't allow small groups to be identified. Even then, there are some 3-digit codes that have fewer than 20,000 residents, and the advice is to lump these together under a new code (000). Looking forward to how GDPR may affect data handling in the UK, is it possible to use a similar approach here?

UK postcodes can be in one of six formats and broken down into seven components that describe geographic areas (Table 1). Given how UK postcodes are constructed, we can't simply truncate postcodes to three characters as is done for USA zip codes as this will result in postcode areas being combined together in ways which do not have real geographic meaning. For example, 'NE3 1ED' and 'NE35 2FG' are both valid postcodes. Truncating both of these to the first three characters would put these both in the 'NE3' group. They should instead both be in an 'NE' area group or separate 'NE3' and 'NE35' district groups (or another chosen postcode component group).

Below, we take a UK postcode population dataset¹ and perform an anonymity analysis of the postcode geographic components on the four population groups present in the dataset: total population, male population, female population and occupied households.

¹http://www.nomisweb.co.uk/output/census/2011/Postcode_Estimates_Table_1.csv

Postcode Format	Outward Code	Inward Code	Area	District	Sub-District	Sector	Unit
AA9A 9AA	AA9A	9AA	AA	AA9	AA9A	AA9A 9	AA
A9A 9AA	A9A	9AA	A	A9	A9A	A9A 9	AA
A9 9AA	A9	9AA	A	A9	N/A	A9 9	AA
A99 9AA	A99	9AA	A	A99	N/A	A9 9	AA
AA9 9AA	AA9	9AA	AA	AA9	N/A	AA9 9	AA
AA99 9AA	AA99	9AA	AA	AA99	N/A	AA99 9	AA

s

Table 1: UK postcode formats and geographic components. 'A' represents an alphabetic character, '9' represent a numeric character. Obtained from <https://ideal-postcodes.co.uk/guides/uk-postcode-format>.

2 Anonymity Analysis

2.1 Anonymity Threshold

For our anonymity analysis we are most concerned about small groups, as they are most likely to contain non-anonymous data. We will use an anonymity threshold, which is the smallest group size that would be acceptable in an anonymised dataset.

The number of groups which are smaller than threshold gives us a measure of how well a chosen combination of postcode component and population group anonymises the data. It is of course possible to aggregate small groups into larger ones, but this is best done when the majority of groups are already anonymised.

2.2 Postcode Component Selection

We have used the area, district, sub-district and sector components of the postcode; these represent increasingly smaller geographic areas, with the full postcode being the smallest. Where a postcode does not have a sub-district (the sub-district is mostly used in areas of London), we have substituted the district in its place.

In Fig. 1 we have plotted the percentage of groups below the anonymity threshold, against the anonymity threshold for each postcode component and population group. Where the percentage of groups below the threshold is low, this means that the data would be well anonymised and vice versa. To use these plots:

1. Choose the population of interest.
 - For example, the total population (top left graph)
2. Select an anonymity threshold on the x-axis.
 - This threshold should be chosen so that data is considered anonymised if all the geographic postcode groups are larger than the threshold
 - i.e. a threshold of 10^5 means that data is considered anonymous if the smallest group size is larger than 100,000.
3. The value of the curve at the chosen threshold tells you what percentage of groups are smaller than this threshold.
 - e.g. for a threshold of 10^5 , only a small number the groups at the postcode area level (red curve) are smaller than 10,000 total population.
 - This would probably be considered anonymous at this level, as the remaining small groups can be aggregated into a single larger one.
 - In contrast, at the sector, sub-district or district level, close to 100% of the groups are below this threshold, and the data would be defined as not anonymous.

We see that each curve has a region where it rises very sharply. To the left this steep region it is clear the data would be well anonymised, and to the right it would be poorly anonymised. Within the steep region it is less clear.

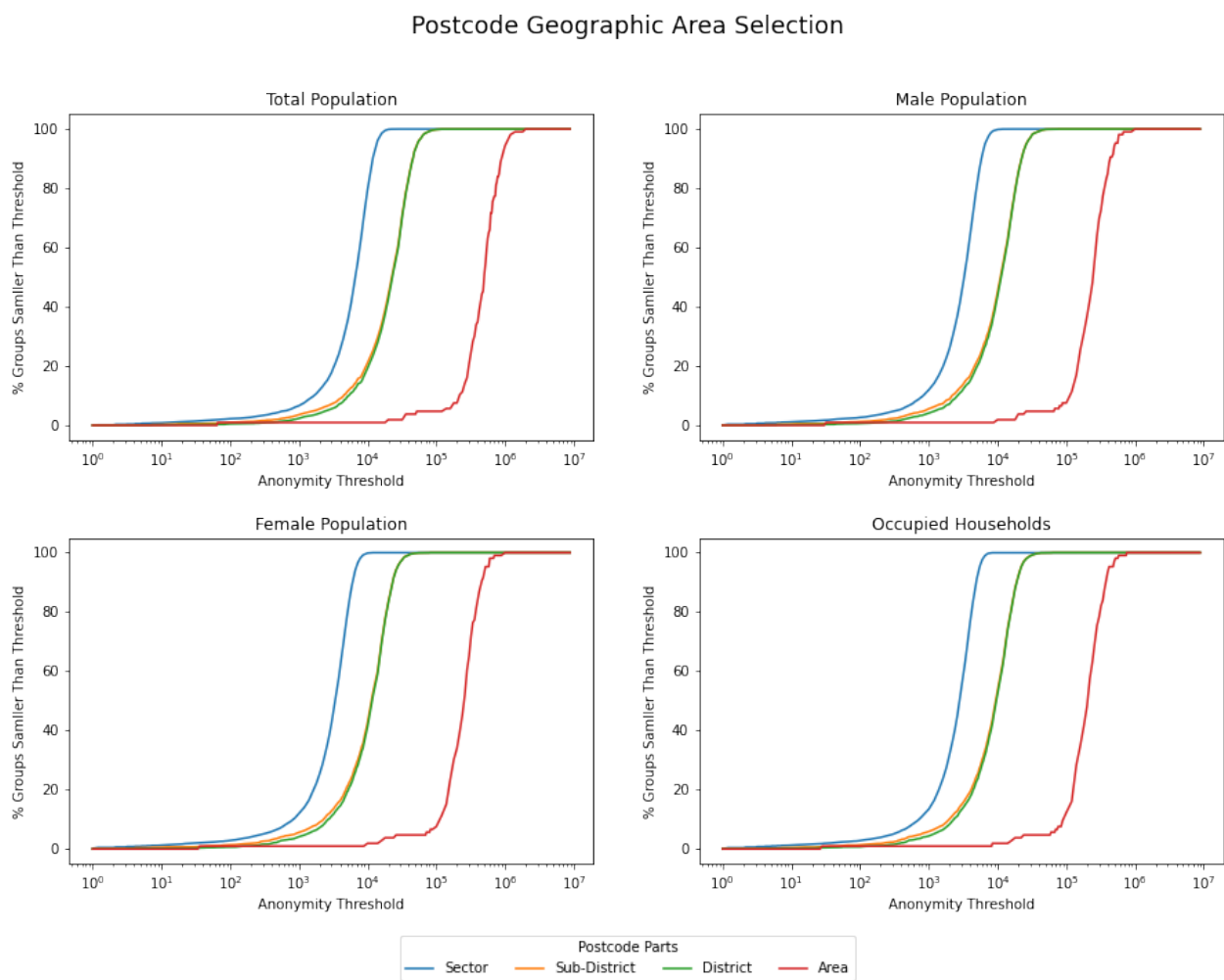


Figure 1: Postcode component selection.

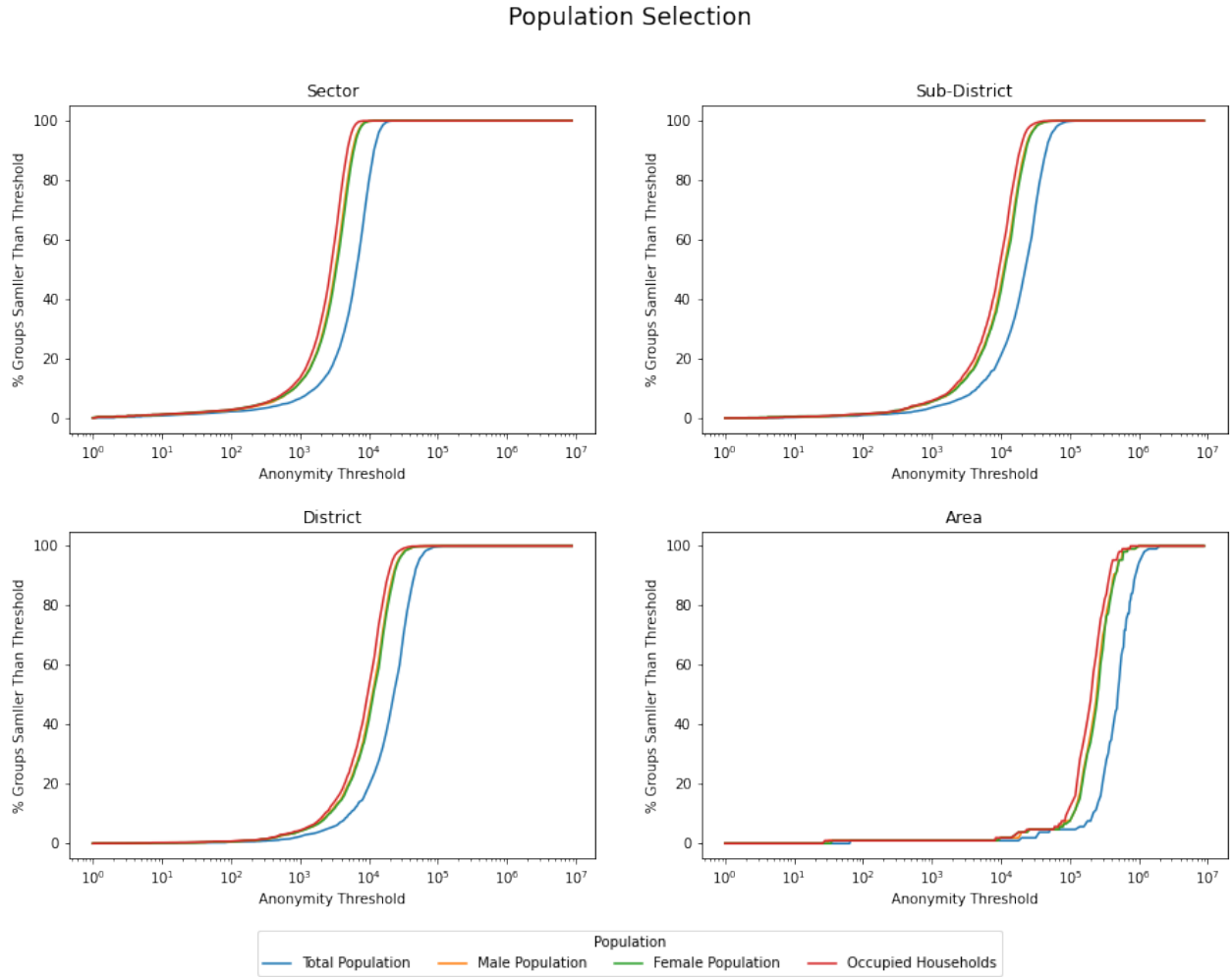


Figure 2: population selection

2.3 Population Segment Selection

If the postcode geographic component is already chosen, it might be desirable to choose what population level to publish: total population, male & female or households. Using the same data we can arrange the plots (Fig. 2) to provide a similar utility to the above for this scenario:

1. Choose the postcode geographic level of interest.
 - For example, the district (bottom left graph)
2. Select an anonymity threshold on the x-axis.
3. The value of the curve at the chosen threshold tells us what percentage of groups are smaller than this threshold.
 - e.g. for a threshold of 10^4 , about 15% of groups at the district level are smaller than 10,000 total population (blue curve) and about 35% of the groups are smaller than the threshold for the remaining population measures.