

Lived population density and the spread of COVID-19

Dave Babbitt* Patrick Garland† Oliver Johnson‡

May 5, 2020

Abstract

We consider variations in the rate of spread of COVID-19, firstly comparing between European countries and secondly comparing between US states. We show that the population density has a small but significant effect on the rate of spread of the virus. However we show that measures of ‘lived population density’, which capture density as perceived by a randomly chosen person, do a better job of explaining variations in the rate of spread, achieving $R^2 = 0.45$ in Europe. We show that adding further measures based on the timing of the outbreak into the regression can increase this to $R^2 = 0.58$.

1 Introduction

Since the onset of the COVID-19 pandemic in December 2019, at least 200,000 people have died of the disease worldwide, with over 3 million positive tests performed [7]. Given the vast human and economic cost, there has been considerable interest [9] in understanding and modelling [3] the spread of the virus, and in comparing casualty figures between and within countries, often through data visualizations.

Such visualizations are often informally used to make comparisons between the responses and interventions of different Governments, or between the quality of their health services. We caution against such judgements, since it is clear that many factors affect the speed of spread and severity of any outbreak, some of which (including demographic factors such as household age profiles, levels of obesity or smoking and cultural factors such as mask wearing or frequency of social contact) exist independently of the quality of any specific Government coronavirus response.

*Data Scientist, Booz Allen Hamilton: email babbitt.dave@bah.com

†School of Medicine, University of Southampton, UK

‡School of Mathematics, University of Bristol, UK: email maotj@bristol.ac.uk

Any assessment of the coronavirus epidemic in any country must be carried out carefully, at the end of the pandemic, in the context of these factors. It will be a major statistical exercise to do this properly, considering interactions rigorously. In this article, we make a minor early contribution to such an audit, by considering the effect of one particular demographic factor, namely population density.

One obvious suggestion is that, since COVID-19 spreads through close proximity between individuals, we might expect faster spread of the virus in countries which have a high population density. We show in Figure 2 that this is true to a reasonable extent, with population density providing an R^2 of 0.23 and p -value less than 0.01 in the two cases studied. This is consistent with previous studies [8], including work [2] modelling an influenza pandemic that led to the model [3], which also indicated a role for population density.

However, we argue that what we will refer to as the *standard population density* of a large region (total number of people divided by area) is not the most useful measure in this context. Indeed, [2] explicitly uses population data on a fine (square kilometre) scale, rather than crude national figures. We will work on similar principles, and consider two versions of what we will call the *lived population density*, which represent the population density as experienced by the average person. We describe these two measures and discuss their properties in Section 2.

We believe that these measures of lived density are more natural to explain the spread of the virus. In Section 3 we explain how we have measured them in practice, as well as the collection of the COVID-19 data. In Section 4 through studying correlations we show that lived density does partially explain variations in the rate of spread of COVID-19 in different regions of a continent. In particular, we show that such measures partially explain variations between European countries and between US states (see Figure 3), achieving an R^2 of 0.45 and 0.39 respectively, with $p < 0.0001$ in each case. In Section 5 we give conclusions, including showing how adding another measure based on timing of the epidemic can increase the R^2 further, and make suggestions for future work.

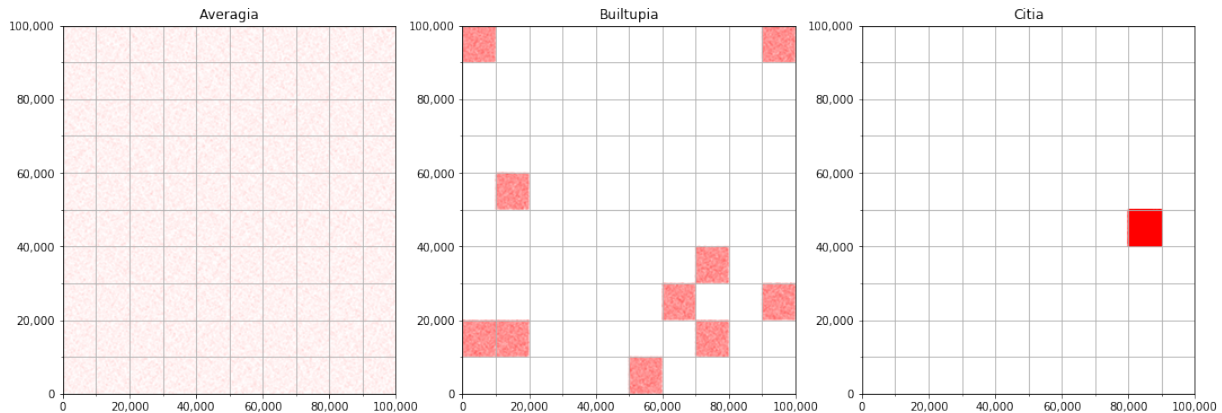


Figure 1: Schematic diagram of population in Averagia, Builtupia and Citia – see Example 2.1.

2 Explanation of different density measures

First we motivate and define the various measures of population density that we use in this paper.

Example 2.1. Consider three countries, each with 100,000 people and an area of 100km², and consider the population of each square kilometre grid square. As illustrated in Figure 1:

1. Averagia has a uniform spread of population, with 1,000 people living in each grid square.
2. Builtupia has ten towns, each consisting of 10,000 people living in a single square kilometre, and the remaining land is uninhabited.
3. Citia has one city, where 100,000 people live in a single square kilometre, and where the remaining land is uninhabited.

Since each standard population density is $100000/100 = 1000$ people/km², we might naively say that each country is equally crowded. However, if we consider the lived experience of a person in each country, it is clear that in everyday life residents of Citia would be close to more people than in Builtupia, who in turn would be close to more people than Averagia. As a result, it seems natural to imagine that COVID-19 would spread more rapidly in Citia than in Builtupia, and in turn faster in Builtupia than in Averagia. We would like to use a measure of population density that captures this, and here define two such measures.

We will consider a region of total area A , divided into M subareas of area A_1, \dots, A_M respectively. We will write n_i for the population of the i th subarea, and write $n = \sum_{i=1}^M n_i$ for the total population of the region.

Definition 2.2. *In this notation, the standard population density is*

$$\rho_S = \frac{n}{A} = \frac{\sum_{i=1}^M n_i}{\sum_{i=1}^M A_i}. \quad (1)$$

Definition 2.3. *We will refer to the non-empty lived density [6] as the same expression normalized by the total area where some people live. That is, we write S for the set of regions with non-zero population, and consider*

$$\rho_N = \frac{\sum_{i \in S} n_i}{\sum_{i \in S} A_i}. \quad (2)$$

Note that in comparison with the expression (1) given for ρ_S , the numerator is unchanged here, since removing terms which are zero does not affect the value of a sum. Indeed, if each subarea contains at least one person then the non-empty lived density ρ_N coincides with the standard population density ρ_S .

Definition 2.4. *We will refer to the quadratic lived density as the sum*

$$\rho_Q = \frac{1}{n} \sum_{i=1}^M \frac{n_i^2}{A_i} = \sum_{i=1}^M \left(\frac{n_i}{n}\right) \left(\frac{n_i}{A_i}\right). \quad (3)$$

We refer to this as a quadratic measure because if each subarea has $A_i = 1$ then it becomes $\rho_Q = \left(\sum_{i=1}^M n_i^2\right)/n$.

Remark 2.5.

1. *We can think of this quadratic lived density ρ_Q as follows: select an individual uniformly at random from the population (with probability n_i/n this will be someone from subarea i). Then n_i/A_i is the density of their local subarea. Hence, assuming the subareas are small enough to be reasonably homogeneous, this measure represents the expected local density, sampling by person.*
2. *In contrast, the standard population density corresponds to sampling by area. That is, pick a point uniformly in the region (with probability A_i/A this will be a point in subarea i). Then according to this distribution, the expected local density will be*

$$\sum_{i=1}^M \left(\frac{A_i}{A}\right) \left(\frac{n_i}{A_i}\right) = \frac{n}{A}, \quad (4)$$

the standard population density ρ_S .

We return to Example 2.1, where recall that each country had standard population density $\rho_S = 1000$. Considering the square kilometre grid squares as the subareas, note that the non-empty lived density ρ_N of Averagia, Builtupia and Citia is $100000/100 = 1000$ people/ km^2 , $10000/10 = 1000$ people/ km^2 and $100000/1 = 100000$ people/ km^2 respectively. In fact, calculation shows that the same three values hold for the quadratic lived density ρ_Q in these cases (though this will not be true in general). For example, for Builtupia, there are 10 non-zero terms in the sum of $\left(\frac{n_i}{n}\right)\left(\frac{n_i}{A_i}\right)$, each equal to $(1/10)(10000/1)$, giving 10000 overall.

Notice that ρ_S , ρ_N and ρ_Q are all measured in the same units, namely people/ km^2 , meaning that it is legitimate to compare them. It turns out that the three measures are always ordered in the same way for every region:

Lemma 2.6. *For any region, the three measures satisfy*

$$\rho_S \leq \rho_N \leq \rho_Q,$$

with equality in the first inequality if and only if all the areas have non-zero population, and with equality in the second inequality if and only if all the areas of non-zero population have the same density.

Proof. The fact that $\rho_S \leq \rho_N$ is clear, since the two expressions both have the same numerator n , but the non-empty lived density has a smaller denominator since the sum is taken over a smaller range. Clearly, equality holds if and only if the denominators are equal.

The fact that $\rho_N \leq \rho_Q$ is more subtle, but can be proved using the Cauchy-Schwarz inequality. Recall that we write S for the set of areas of non-zero population, then use the fact that

$$n^2 = \left(\sum_{i \in S} \frac{n_i}{\sqrt{A_i}} \sqrt{A_i} \right)^2 \leq \left(\sum_{i \in S} \frac{n_i^2}{A_i} \right) \left(\sum_{i \in S} A_i \right) = (n\rho_Q) \left(\sum_{i \in S} A_i \right), \quad (5)$$

and rearrange, using the fact that the value of ρ_Q is unchanged if we restrict to a sum over $i \in S$ (since terms with $n_i = 0$ do not contribute to the sum). Equality holds if and only if each populated subarea has the same density, that is if n_i/A_i does not depend on i for each $i \in S$. \square

Indeed, a version of the same argument shows that if we partition a subarea into smaller subareas (by obtaining more finely grained population and area data) then the quadratic lived density ρ_Q does not decrease, and will strictly increase unless all the new subareas have the same density. In other words, the maximum value of ρ_Q for a region is obtained by breaking it into small homogeneous subareas.

Definition 2.7. For a given region we define the regularity coefficient $R = \rho_Q/\rho_S$. Note that this is a dimension-free quantity.

From Lemma 2.6, we know that $R \geq 1$ for any region. This ratio measures the extent to which population is evenly distributed in a region; returning to the toy Example 2.1, notice that Averagia has $R = 1$, Builtupia has $R = 10$ and Citia has $R = 100$. These values may help calibrate our understanding of the values of R observed for actual regions.

We believe that lived density measures quantities merit further investigation. For example, instead of simply considering quadratic lived density measures, it may be that local densities have a larger effect, and for some $p > 1$ that a p -norm type quantity of the form

$$\rho_Q^{(p)} = \left[\sum_{i=1}^M \left(\frac{n_i}{n} \right) \left(\frac{n_i}{A_i} \right)^p \right]^{1/p} \quad (6)$$

may explain better the variation in the rate of spread of COVID-19. Such quantities would put even greater weight on areas of high density, which may be appropriate in terms of the effect on the spread.

Clearly the ρ_Q of Definition 2.4 corresponds to the case $p = 1$ in the sense of (6). More curiously, the ρ_S of Definition 2.2 corresponds to taking $p = -1$ in (6). In other words, ρ_Q is a weighted arithmetic mean of the local densities and ρ_S is the corresponding weighted harmonic mean. In that sense, the fact that (see Lemma 2.6) $\rho_S \leq \rho_Q$ can be seen as a consequence of the weighted arithmetic mean–harmonic mean inequality, and in general the expressions $\rho_Q^{(p)}$ from Equation (6) will be increasing in p – see [4, Eq. (2.9.1)].

3 Data

We first describe the data used to model the effect of these different population density measures on the spread of COVID-19.

3.1 Death data

Since COVID-19 epidemics have currently reached different stages in different countries, we do not focus on the total number of deaths that have now occurred, but rather on the rate of spread early on. In the sense of the classical SIR model [5], we anticipate that this is a phase of unrestricted exponential growth, where we seek to compare the growth exponent.

We do this by choosing a relatively small threshold value of deaths, and then comparing the number of deaths that have occurred a fixed number of days later. For concreteness,

we choose these numbers to be 5 and 10: that is, our metric is “number of deaths that have occurred 10 days after deaths hit 5”.

Clearly the choice of these numbers is somewhat arbitrary: we want to choose a threshold that is large enough that daily numbers are not affected too much by random fluctuations, and to wait a long enough period for random daily effects to cancel out. However, waiting too long means that the numbers can be affected by Government actions such as lockdown. Further, waiting for deaths to hit too high a level may exclude some smaller countries. Some experimentation has shown that our results are robust to the choice of parameters.

Although there is clearly variation between countries in terms of how deaths are recorded and announced, we hope that considering spread in this way (rather than through absolute values) will reduce such effects. Essentially, since the epidemic is growing exponentially [5], we hope that systematic under-counting or lags in reports will not affect the estimate of the growth rate.

The data on COVID-19 deaths itself was obtained online. Data for European countries was taken from [7], which is itself based on ECDC data. Data for US states was downloaded from [1].

3.2 Population data

We in fact use two different measures to make comparisons within the two continents. This is largely an issue of availability of data; clearly it would be preferable to use both measures in both cases, and compare the results, to understand the applicability of such methods.

We studied 30 European countries, using figures of standard population density ρ_S taken from Wikipedia, and using the values of the non-empty lived density ρ_N for each country calculated by Rae from Eurostat data, and stated in [6]. Some extremely small countries (where the size of the COVID-19 outbreak was too small to achieve the threshold discussed above) were excluded from the comparison.

We studied the 50 US States and District of Columbia, using figures of standard population density ρ_S taken from Wikipedia, and calculating the quadratic lived density ρ_Q ourselves. The methodology for this calculation involved breaking the United States into 3,142 subareas (at the level of county, parish, borough or city). Population values were available for each such subarea from the US census website [10], and the Wikipedia API was used to find the land area for each region. Combining these values as described in Definition 2.4 allowed us to find the value of ρ_Q for each state.

We tabulate these values in Appendix A for completeness. Note the wide variation in quadratic lived density values ρ_Q , ranging from $\rho_Q = 5$ for Wyoming up to $\rho_Q = 9376$ for Connecticut. It may also be interesting to observe the range of values of regularity

coefficient R and its relation to the overall density figures. Low values of R include 1.1 for the District of Columbia (both ρ_S and ρ_Q very high) and 1.8 for Vermont (both ρ_S and ρ_Q very low), which indicate a uniform spread of population. The highest value of R was 39.4 for New York state (ρ_S fairly low, ρ_Q very high) indicating a very uneven spread of population.

4 Results

In Figure 2 we plot correlations between the rate of spread and standard population density, first for Europe and then for the US. Both axes are plotted on a log scale, since in the context of an SIR epidemic [5] we would like to measure the growth exponent. In particular we plot the logarithms of the density values for US states listed in Appendix A. Observe that in each case there is a significant correlation: the p -values of the regression coefficient corresponding to slope are $p = 0.007$ and $p = 0.0005$ respectively, indicating a statistically significant dependence on ρ_S .

However in both cases, the R^2 of 0.23 is relatively low, and it is natural to wonder whether we can do better. An examination of certain outlying points gives a hint as to an issue here. The European country exhibiting the highest rate of spread was Spain, with 309 deaths, but a relatively low standard population density $\rho_S = 93$ people/ km^2 . However, as discussed in detail in [6], much of Spain is unpopulated, with very high density population in Barcelona and Madrid. As a result, other than three very small countries (Monaco, Andorra and Malta, not studied here because of the size of their epidemic), Spain has the highest lived non-empty population density in Europe, with a value of $\rho_N = 737$ people/ km^2 .

A similar picture emerges in the US data. The highest rate of spread was observed in New York state, with 385 deaths, but again a relatively low standard population density $\rho_S = 162$ people/ km^2 . Again, this relatively low value of ρ_S can be attributed to large low population areas outside New York City. When using the lived measure, New York state in fact has the second largest value of $\rho_N = 6386$ people/ km^2 , beaten only by Connecticut. In the context of Definition 2.7 this indicates a very high regularity coefficient of 39.4, not far off that of the toy example Citia.

As a result of such effects, when we plot correlations between the rate of spread and the respective lived density measures, the effect is stronger than for the corresponding plots for ρ_S . In Figure 3a) we see a significant correlation between rate of spread and ρ_N for Europe ($R^2 = 0.45$ and $p < 0.0001$). In Figure 3b) we see a similar result for the United States in terms of ρ_Q ($R^2 = 0.39$ and $p < 0.0001$).

In other words, we conclude that in each case using non-standard measures of population density reveals more statistically significant effects than those arising from the standard population density ρ_S .

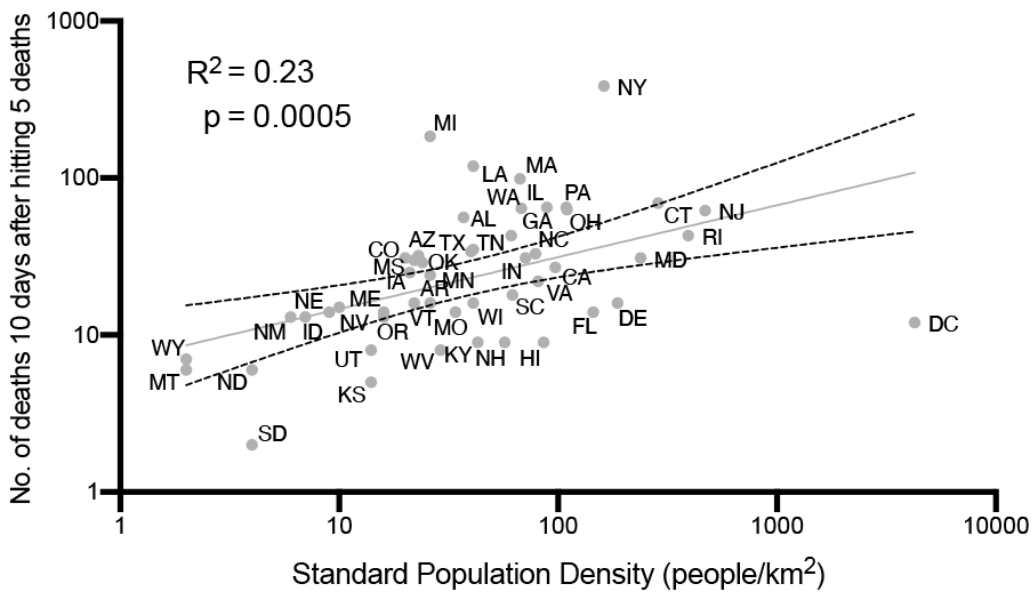
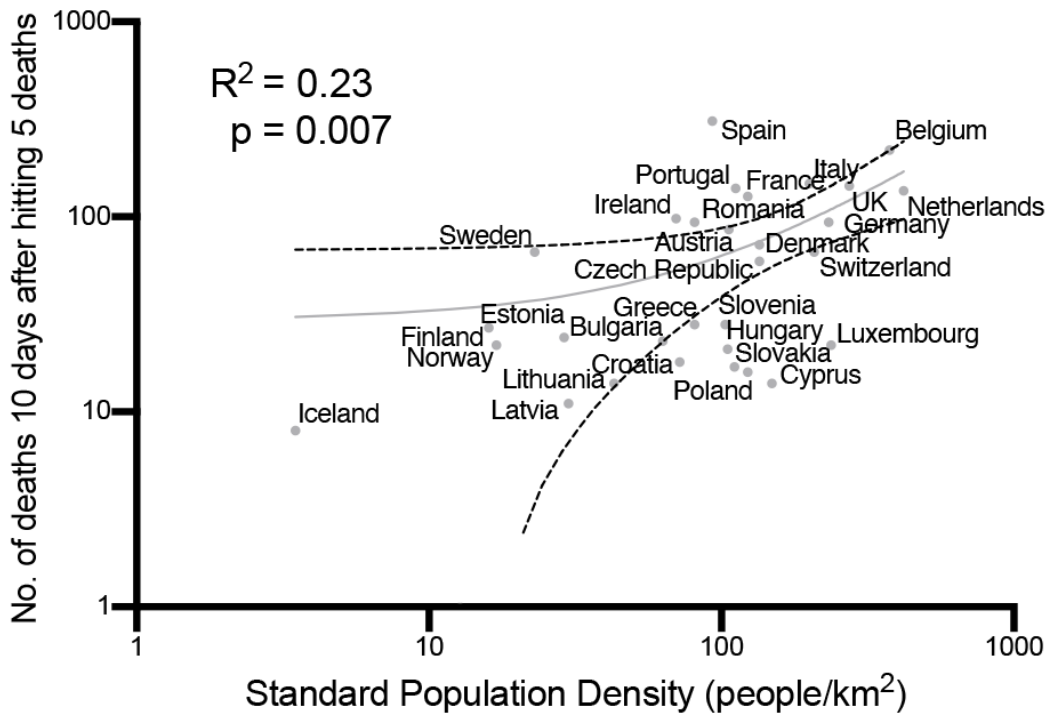


Figure 2: Scatterplot of rate of spread against standard population density ρ_S : a) for European countries b) for US states. Both plots are on a log-log scale

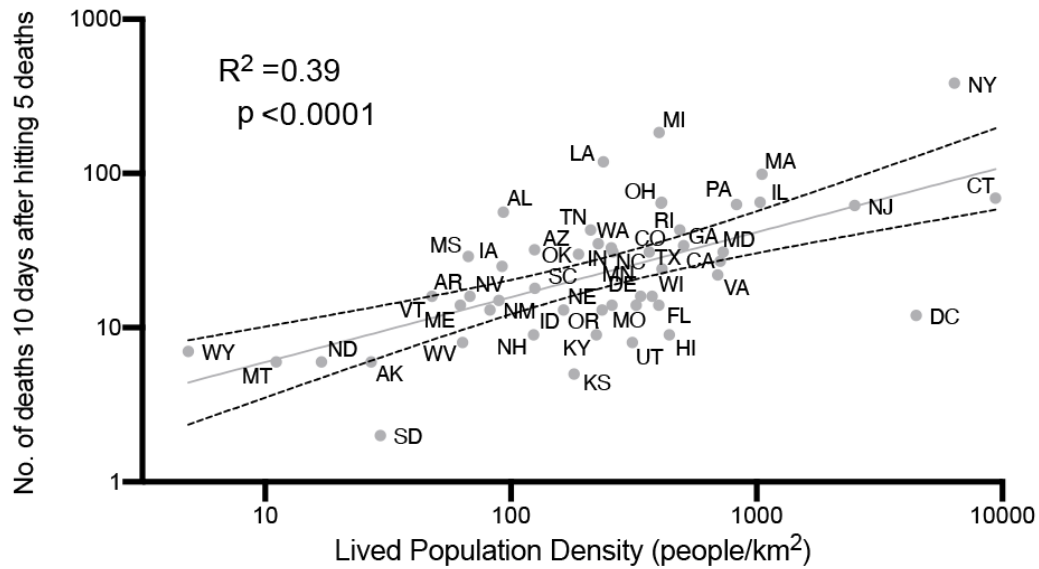
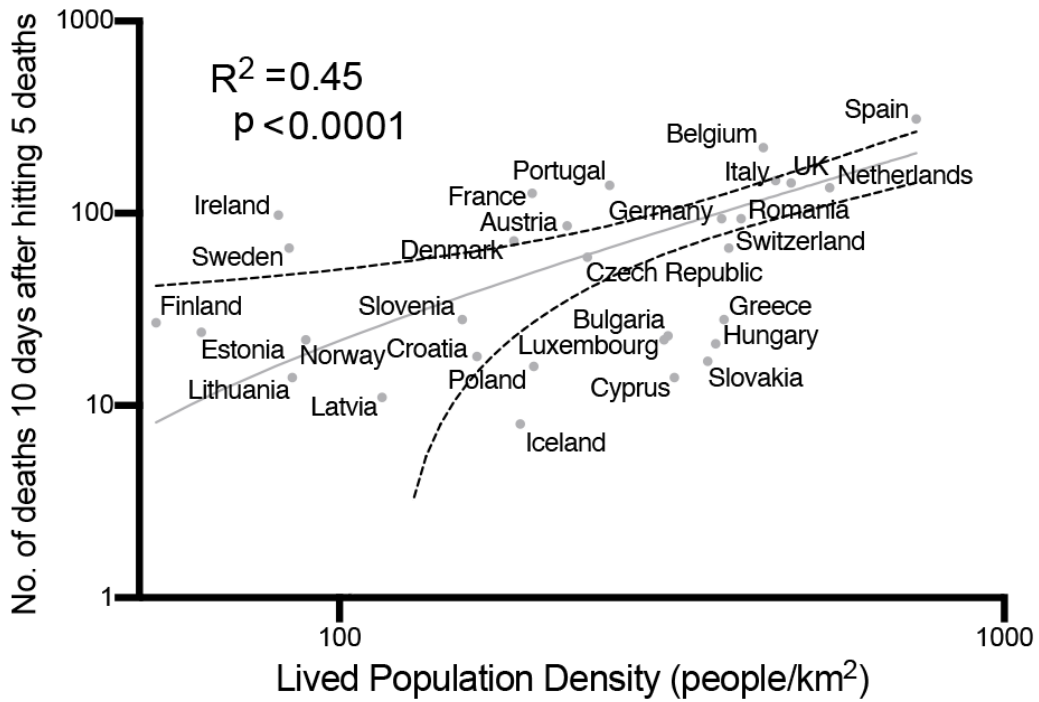


Figure 3: Scatterplot of nonempty lived population density against rate of spread: a) plots rate of spread against ρ_N for European countries b) plots rate of spread against ρ_Q for US states. Both plots are on a log-log scale

5 Discussion and further results

We have discussed properties of two non-standard measures of population density ρ_N and ρ_Q and argued by examination of correlation plots that these quantities explain more of the variation in the rate of spread of COVID-19 than the standard population density ρ_S .

While we still caution against making premature comparisons on the basis of COVID-19 numbers, considering countries which lie above or below the confidence regions in plots such as Figure 3 may provide a more principled way to do this than relying on raw numbers alone. Such comparisons may suggest other explanatory factors that should be added into a model, or may prompt further examination of Government policies or interventions that have led to unsuccessful or successful management of the epidemic.

In terms of epidemic management, understanding the role of ρ_N or ρ_Q may suggest a differential approach to, for example, test and trace or social distancing. So, test and trace could be focused on high lived density areas, and social distancing might be less strict in low lived density areas (which could allow those areas to contribute optimally to economic activity). Similarly it may allow screening of international travellers to focus on arrivals from high lived density countries early in an epidemic.

One drawback of such methods is the amount of work required to calculate the non-standard measures ρ_N and ρ_Q for given regions. As mentioned previously, ideally ρ_Q would be calculated on the basis of a division into subareas of homogeneous density. In some cases population and area data may not be available on a sufficiently fine scale, and clearly the finer the scale the more work is required to compile it. However, we note that, once calculated, ρ_Q and ρ_N can be reused for a variety of modelling purposes, and so their evaluation can be regarded as ‘once-only work’, and we provide the US values in Appendix A for reference.

One obvious question is whether this method allows for comparisons between every country in the world, and whether the local density measures universally predict the spread of COVID-19. Unfortunately this appears not to be the case, despite the success demonstrated above on a continental scale. For example, Hong Kong has a standard population density of $\rho_S = 6659$ people/ km^2 (with its ρ_Q presumably higher still). Even this ρ_S value is already higher than the $\rho_Q = 6385$ people/ km^2 for New York state, and yet Hong Kong has had an extremely low number of deaths so far, in common with many other Far East countries [7].

We believe that this may be because of other significant factors, principally cultural norms around the widespread wearing of masks, with the possibility that factors such as climate may also play a role. However, we believe that a (correctly defined) measure of population density may play an important role in understanding the rate of spread, and should feature in the discussion when making comparisons of the final outcome of the pandemic within continents.

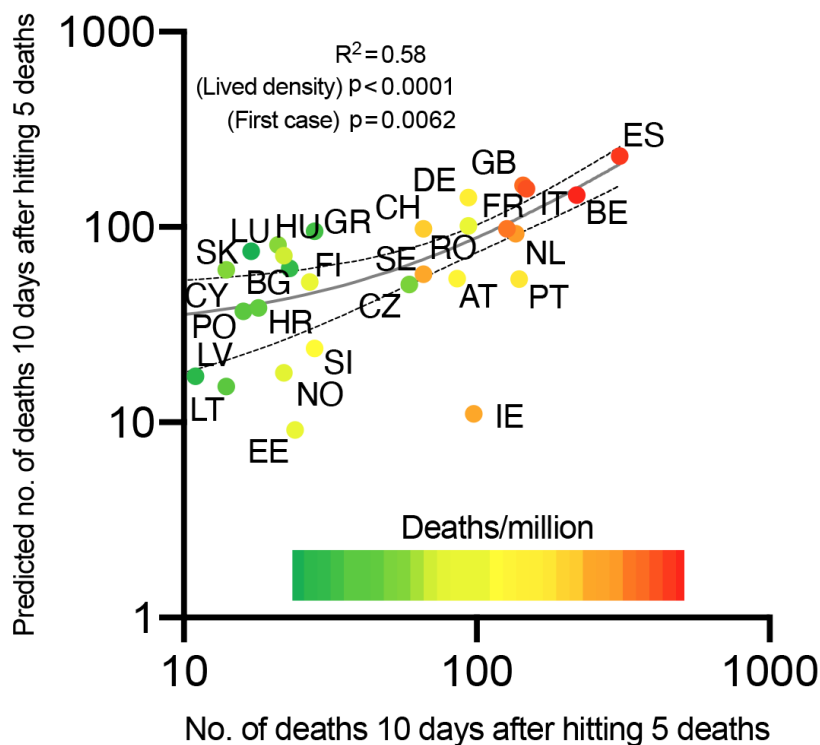


Figure 4: Scatterplot showing how multiple linear regression, taking into account lived population density ρ_N and first case time, can explain much of the variation in rate of spread.

Using aggregate data derived from over 1 billion people between two continents, we have shown that an appropriately defined measure of population density can capture almost half of the variation in the rate of COVID-19 fatalities. What other variables might also account for the rate of spread within Western countries?

Already it is emerging that COVID-19 morbidity and mortality is influenced by demographics (eg, age, sex), obesity, and pre-existing health conditions (diabetes, hypertension). However, we might also hypothesize that the rate of spread could be affected by how frequently outbreaks were seeded in countries due to international travel, or how much time each country had to prepare for such outbreaks after the seriousness of COVID-19 was known. The latter of these two could relate to government action, but also simple changes of behaviour such as recognising that cold and flu symptoms require social distancing and isolation.

To capture these possible factors we calculated the time between when the World Health Organisation declared COVID-19 could transmit between people (20th January 2020) to when each European country had its first clinically confirmed case. As an example, Poland

had 44 days before its first confirmed COVID19 case, whereas the U.K. had only 11 days.

Multiple linear regression was used to examine whether lived density and time to ‘first case’ could independently explain variation in the rate of COVID-19 spread in European countries (Figure 4). These two variables were able to explain 0.58 of the variation within the rate of COVID-19 infection (lived density $p < 0.0001$; ‘first case’ $p = 0.0062$).

The relationship between the initial rate of spread of COVID-19 between European countries and their subsequent burden of disease is also illustrated by indicating how many death per million each country currently (as of 2nd May 2020) has – see colouring of points in Figure 4. It is clear that the worst affected countries per capita so far have been those where the early spread was fastest, and so modelling and understanding the initial spread is useful in this context.

References

- [1] Covid Tracking project. Covid tracking website. Available at <https://covidtracking.com/api/v1/states/daily.csv>, Downloaded 26th April 2020, 2020.
- [2] N. Ferguson, D. Cummings, C. Fraser, J. C. Cajka, P. C. Cooley, and D. S. Burke. Strategies for mitigating an influenza pandemic. *Nature*, 442(7101):448–452, 2006.
- [3] N. Ferguson, D. Laydon, G. Nedjati Gilani, N. Imai, K. Ainslie, M. Baguelin, S. Bhatia, A. Boonyasiri, Z. Cucunuba Perez, G. Cuomo-Dannenburg, et al. Report 9: Impact of non-pharmaceutical interventions (NPIs) to reduce COVID19 mortality and healthcare demand. 16th March 2020. doi:10.25561/77482.
- [4] G. H. Hardy, J. E. Littlewood, and G. Pólya. *Inequalities*. Cambridge University Press, 1934.
- [5] W. O. Kermack and A. G. McKendrick. A contribution to the mathematical theory of epidemics. *Proceedings of the Royal Society of London. Series A*, 115(772):700–721, 1927.
- [6] A. Rae. There’s a better way to measure population density. Available at <https://www.citylab.com/life/2018/02/theres-a-better-way-to-measure-population-density/552815/>, Downloaded 26th April 2020, 8th February 2018.
- [7] M. Roser, H. Ritchie, E. Ortiz-Ospina, and J. Hasell. Our World In Data: Coronavirus Disease (COVID-19). Available at <https://ourworldindata.org/coronavirus>, Downloaded 2nd May 2020, 2020.
- [8] P. M. Tarwater and C. F. Martin. Effects of population density on the spread of disease. *Complexity*, 6(6):29–36, 2001.

- [9] C. Tobitt and W. Turvill. Times, FT and WaPo discover 'real appetite' for data-driven visual journalism on coronavirus. *Press Gazette*, 9th April 2020.
- [10] US Census Office. United States census data. Available at <https://www2.census.gov/programs-surveys/popest/tables/2010-2019/counties/totals/co-est2019-annres.xlsx>, Downloaded 26th April 2020, 2020.

A Table of densities for US States

State	Lived ρ_Q	Standard ρ_S	Regularity R
AL	93	37	2.5
AK	27	1	27.0
AZ	125	23	5.4
AR	68	22	3.1
CA	711	97	7.3
CO	366	20	18.3
CT	9376	286	32.8
DE	337	187	1.8
DC	4464	4251	1.1
FL	400	145	2.8
GA	408	68	6.0
HI	442	86	5.1
ID	164	7	23.4
IL	1032	89	11.6
IN	260	71	3.7
IA	92	21	4.4
KS	181	14	12.9
KY	222	43	5.2
LA	238	41	5.8
ME	62	16	3.9
MD	732	238	3.1
MA	1053	336	3.1
MI	402	67	6.0
MN	412	26	15.8
MS	67	24	2.8
MO	324	34	9.5
MT	11	2	5.6
NE	258	9	28.7
NV	89	10	8.9
NH	124	57	2.2
NJ	2505	470	5.3
NM	82	6	13.7
NY	6386	162	39.4
NC	257	79	3.2
ND	17	4	4.2

State	Lived ρ_Q	Standard ρ_S	Regularity R
OH	411	109	3.8
OK	189	22	8.6
OR	235	16	14.7
PA	828	110	7.5
RI	487	394	1.2
SC	125	62	2.0
SD	29	4	7.4
TN	211	61	3.5
TX	505	40	12.6
UT	312	14	22.3
VT	48	26	1.8
VA	695	81	8.6
WA	227	41	5.5
WV	64	29	2.2
WI	376	41	9.2
WY	5	2	2.4