

# Generalized Radiation Model for Human Migration

Christian M. Alis (✉ [calis@aim.edu](mailto:calis@aim.edu))

Asian Institute of Management

**Erika Fille Legara**

Asian Institute of Management

**Christopher Monterola**

Asian Institute of Management

---

## Research Article

**Keywords:** human migration, radiation model , outperforms , number of amenities

**Posted Date:** March 22nd, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-319100/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published at Scientific Reports on November 22nd, 2021.  
See the published version at <https://doi.org/10.1038/s41598-021-02109-1>.

# Generalized Radiation Model for Human Migration

Christian M. Alis, Erika Fille Legara, and Christopher Monterola

Analytics, Computing and Complex Systems Laboratory (ACCeSs@AIM), Asian Institute of Management, 123 Paseo De Roxas, Makati City, 1229 Philippines

## ABSTRACT

One of the main problems in the study of human migration is predicting how many people will migrate from one place to another. An important model used for this problem is the radiation model for human migration, which models locations as attractors whose attractiveness is moderated by distance as well as attractiveness of neighboring locations. In the model, the measure used for attractiveness is population which is a proxy for economic opportunities and jobs. However, this may not be valid, for example, in developing countries, and fails to take into account people migrating for non-economic reasons such as quality of life. Here, we extend the radiation model to include the number of amenities (offices, schools, leisure places, etc.) as features aside from population. We find that the generalized radiation model outperforms the radiation model by as much as 7.7% relative improvement in mean absolute percentage error based on actual census data five years apart. The best performing model does not even include population information which suggests that amenities already include the information that we get from population. The generalized radiation model provides a measure of feature importance thus presenting another avenue for investigating the effect of amenities on human migration.

## Introduction

Understanding and predicting the rate of flow between locations have applications in urban and transport planning<sup>1,2</sup>, epidemic modelling<sup>3-6</sup> and emergency management<sup>7,8</sup>, among others. For many years, the gravity model and its variations<sup>9</sup> have been the go-to model for predicting these movements. In this model, migration flow is proportional to the population of the source and destination localities, and inversely proportional to their distance.

More recently, the radiation model (RM) for human migration<sup>10</sup> was introduced and predicts the average flow of migrants  $\langle T_{ij} \rangle$  from locality  $i$  to locality  $j$  as

$$\langle T_{ij} \rangle = T_i \frac{p_i p_j}{(p_i + s_{ij})(p_i + p_j + s_{ij})} \quad (1)$$

where  $T_i$  is the total number of migrants from  $i$ ,  $p_i$  and  $p_j$  are the population in  $i$  and  $j$ , respectively, and  $s_{ij}$  is the total population in the circle centered at  $i$  and touching  $j$  excluding the source and the destination populations. It has been shown that this model and its variations can replicate the observed changes in population across several cities in developed countries<sup>2,10-13</sup> but less so in developing countries<sup>6,14</sup>.

The idea behind the model is that migrants are motivated to move towards localities with better economic opportunities such as availability of jobs. However, the pull of one locality is tempered by the pull of neighboring localities as well: a highly urbanized city would have a stronger pull if it is surrounded by rural areas compared to it being part of a metropolis. Similarly, the model gives preference to migration between localities that are nearer to each other over longer distance migrations.

Instead of using actual economic indicators to measure the economic opportunities in a locality, the model instead used population as a proxy: the bigger the population of a locality the more economic opportunities it has. However, for developing countries, this assumption may not hold. Due to higher likelihood of inequality in a developing country, a bigger population may not necessarily imply more economic opportunities. In fact, because of increased competition for limited economic opportunities in a crowded city, residents may be tempted to move out to less crowded localities with relatively more opportunities per capita. Moreover, even if the economic opportunities per capita is better, poorly regulated cities in developing countries are challenged by lower quality of life due to crimes, pollution, traffic congestion, and weak peer/community support system. For developing countries, religion or faith-related culture can also play a role in one's migration decision (for example, Muslims are not welcome in some Christian-majority areas and vice versa) as discrimination is often enhanced by poor quality of education.<sup>15,16</sup> The importance of tribal, cultural and linguistic differences has already been shown to affect human mobility significantly more than that for a developed country<sup>5,14</sup>.

Even in developed countries, individuals may want to move to a locality in search for a better quality of life instead of better employment<sup>17</sup>. Some migrants do not stay in one city and sometimes even return to where they were before.<sup>18</sup> Indeed, it is already well known that some residents in urban areas opt to move to rural areas following a process collectively known

as counterurbanization<sup>19</sup>. More recently, evidence has been found that the lateral movement of people from one rural area to another is a significant chunk of rural in-migration<sup>20,21</sup> hence we cannot always assume that people from rural areas will move to urban areas if ever they move.

In this paper, we propose a generalized radiation model (GRM) for human migration. Instead of using population as the only proxy, we combine it with other characteristics of the locality to form an urbanization index  $U$ . We then use  $U$  for estimating  $\langle T_{ij} \rangle$  instead of population:

$$\langle T_{ij} \rangle = T_i \frac{U_i U_j}{(U_i + v_{ij})(U_j + v_{ij})} \quad (2)$$

where  $U_i$  and  $U_j$  are the urbanization index at  $i$  and  $j$ , respectively, and  $v_{ij}$  is the total urbanization index in the circle centered at  $i$  and touching  $j$  excluding the source and the destination population.

Aside from improved applicability of GRM to more countries, another benefit of GRM is that it provides us another method for investigating the drivers of migration since  $U$  is composed of several components.

There were already attempts at generalizing the radiation model. Kang et al.<sup>13</sup> introduced a correction for spatial scale as well as the amount of push from the source, pull towards the destination, and interventions in between. Liu and Yan introduced the opportunity priority selection<sup>22</sup> and universal opportunity<sup>23</sup> models that generalize how trip selection by individuals is influenced by the opportunities at destinations and intervening opportunities from source to destination. None of these generalized models, however, directly model how multiple features such as amenities contribute to the attractiveness of opportunities of a place or locality.

Similar to Robinson and Dilkina<sup>24</sup> which directly estimated  $\langle T_{ij} \rangle$  from exogenous data and to McCulloch et al.<sup>25</sup> which created an ensemble of models, we used machine learning to build the model. However, by anchoring GRM on the radiation model, our model is more mechanistic and easier to interpret compared to a pure machine learning model. By using the amenities in each locality as our feature, we are able to estimate the attractiveness of each amenity type for migrants.

In the next section, we elaborate more on GRM especially on how to estimate  $U$  and  $T_i$ . We then explore the fitted GRM for a developing country, the Philippines, and compare it with that of RM.

## Urbanization index

The urbanization index  $U$  is the analog of population in GRM. It is simply the weighted sum of component factors  $f_k$ ,

$$U = \sum_k w_k f_k, \quad (3)$$

where  $w_k$  is the weight of the  $k$ th factor. The factors  $f_k$  can be any feature of a locality. In this paper, we consider the population, population density and the number of structures in the locality for different kinds of amenities as the features.

The intuition behind the use of amenities as features is that a migrant may want to move to a locality based on what is important to them according to the culture of the segment of population which they belong. While the default feature of job opportunities can be represented by the number of offices and/or office space it also weights independently other factors. For example, presence of schools might be attractive for a family hoping to have the next generation lift them out of poverty. Another case, for relatively well-off families, the presence of considerable leisure places might be more attractive as it suggest a better work-life balance. These various motivations can be readily modeled by the use of amenity counts in a locality, resulting in a more granular and less biased estimation of quality of life. For the standard radiation model, all of such scores are generically simplified to be based on the relative population density. Indeed, the importance of amenities or places were already hinted empirically by Noulas et al.<sup>26</sup> and is a central concept of Stouffer's intervening opportunities theory<sup>27</sup>.

The  $f_k$ 's have varying scales hence the values of each feature should be normalized to make the features comparable to each other. Instead of picking one, we investigate three methods of normalization:

- *Min-max*: The value of each feature is rescaled to  $[0,1]$  corresponding to the minimum and maximum values of that feature,

$$x \rightarrow (x - x_{\min}) / (x_{\max} - x_{\min}). \quad (4)$$

This normalization implies that the magnitude of a feature matters. Thus, if the locality with the maximum value is an outlier then that locality would yield a much stronger pull for migrants while the other localities would have similar pull for that feature.

- *Adjusted z-score*: The z-score of each feature value is computed but since the transformed value cannot be negative, we translate the value to the right by one standard deviation then set to zero those that are still negative,

$$x_{\text{adj}} = (x - \bar{x})/\sigma + \sigma \quad (5)$$

$$x \rightarrow \begin{cases} x_{\text{adj}} & x_{\text{adj}} \geq 0 \\ 0 & x_{\text{adj}} < 0 \end{cases}. \quad (6)$$

This normalization also implies that the magnitude of a feature matters. Compared to *Min-max*, there is a stronger bias towards localities that have high values for that feature because those that have low values (more than 1 standard deviation less from the mean) would have zero weights while those with positively outlying values would be more emphasized.

- *Percentile*: The percentile of the value for that feature is used. This implies that migrants only look at the relative rank of the locality and not on the actual value for that amenity. Thus, outliers would not distort the implied pull for that amenity.

Together with the estimated weights of each feature, the above normalization procedure provides an anchor for interpreting the hierarchy and dynamics of the features used with respect to model accuracy. This will be highlighted in the discussion of results.

Both  $T_i$  and  $w_k$  can be considered as trainable parameters. The change in population  $\Delta p_i$  of locality  $i$  is

$$\Delta p_i = p_i^{t+1} - p_i^t = (b - d) p_i^t + \sum_{r \neq i} \langle T_{ri} \rangle - \sum_{r \neq i} \langle T_{ir} \rangle, \quad (7)$$

where  $p_i^t$  and  $p_i^{t+1}$  are the population at times  $t$  and  $t + 1$  for locality  $i$ , and  $b$  and  $d$  are the birth rate and death rate, respectively. To train  $T_i$  and  $w_k$ , we need the population of the localities for two time points to compute  $\Delta p_i$ . We can then use stochastic gradient descent to minimize the mean squared error (MSE) between the observed  $\Delta p_i$  and the estimated  $\Delta p_i$ .

The trained  $w_k$  and  $U$  are highly interpretable. A more positive value of  $w_k$  implies that  $f_k$  drives people to move towards the locality. Similarly, a more positive  $U$  has more pulling power compared to other neighboring localities.

## GRM in a developing country

We compare the results of GRM using three normalization methods and classic RM for a developing country, the Philippines. We also add a baseline model wherein we scale the local change in population according to the national change in population according to the census. This model implies a uniform birth, death and migration rate, which are equal to the national rates, in all localities.

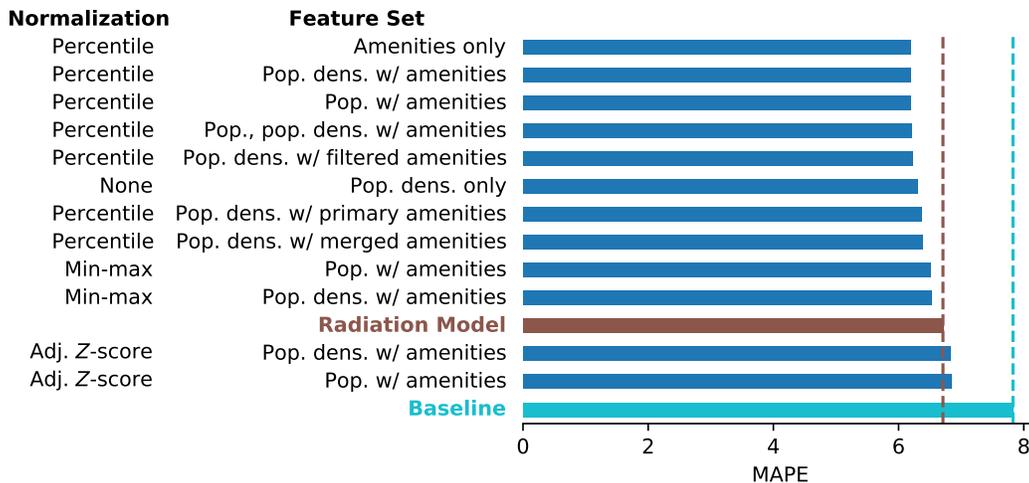
Census data of the Philippines was taken from the Philippine Statistics Authority website ([psa.gov.ph](http://psa.gov.ph)). The three most recent censuses were conducted in 2007, 2010 and 2015, thus, we use the 2007 census as the base year, 2010 census for calibrating the model and 2015 census as the test year. We consider the administrative level 2 (city and municipality) population, which we simply refer to as locality.

Forecasting the population for a year that is after the calibration year (2010) is done by iteratively creating annual forecasts until the desired year is reached. We start by forecasting the 2011 population (calibration year + 1) based on the projected amenity counts for 2010, and the birth and death rates for 2010. Since the trained model is calibrated for a three-year timestep (2007-2010), the raw prediction is divided by three to make it annual. Based on the forecasted population for 2011, we project the amenity counts for 2011, which we then use to forecast the population for 2012 (calibration year + 2) and so on. In all cases, we use the actual birth rate and death rate of the previous rate. This information is, of course, not available if we are really forecasting five years into the future but since we are only using the forecasts to compare models, the use of actual rates should be acceptable. The quoted performance metric values in this paper should therefore be considered to be the best possible values of the models.

For comparing the performance of models, we look at the mean absolute percentage error (MAPE) between the forecasted locality population in 2015 with the census population,

$$MAPE = \frac{1}{N} \sum_i \left| \frac{\hat{p}_i^{2015} - p_i^{2015}}{p_i^{2015}} \right| \times 100\%, \quad (8)$$

where  $N$  is the number of localities,  $\hat{p}_i^{2015}$  is the forecasted population in 2015 for locality  $i$  and  $p_i^{2015}$  is the population in 2015 for  $i$  according to the census. This measure is more relevant than the usual MSE because the distribution of population across localities is fat-tailed so MSE will be heavily biased towards localities with larger population.



**Figure 1.** Model performance based on mean absolute percentage error (MAPE). Percentile normalization is the best method of normalization, outperforming Min-max and Adjusted Z-score normalizations in all instances. The best performing model uses only amenities as features and follows Percentile normalization. It corresponds to a 7.7% MAPE improvement relative to Radiation Model. All models beat the baseline model which is the outright scaling of locality population according to the same rate of change in the national population. The performance metrics of all of the 78 configurations that were investigated are displayed in Supplementary Table S1.

We also considered minimizing the forecast MAPE directly as well as minimizing the MSE of the difference in the logarithmic forecast and logarithmic actual population. However, both resulted in worse performance so they are no longer further described in this paper.

For every model configuration, 100 realizations of the model are trained. We considered taking the mean weight of a feature as well as the median weight of the feature when doing the forecast. However, results show that taking the mean clearly outperforms taking the median so for the rest of this paper, we only quote the results for models that took the mean feature weight.

The complete table of performance metrics for all the 78 configurations that we investigated is in Supplementary Material Table 1.

## Backcasting the number of amenities

We use OpenStreetMap (OSM) data for counting the number of amenities per locality. Administrative boundaries are courtesy of GADM v3 (gadm.org). Amenity information is based on the available information on OSM on 1 Aug 2015, the first day of the 2015 census, and reconstructed from the 24 Feb 2020 historical OSM data dump.

Although the OSM road coverage for the Philippines is quite high<sup>28</sup>, we are not as confident with points-of-interest (POI) coverage especially in the earlier years 2007 and 2010, corresponding to the base and target census years. To minimize issues of coverage, we instead create a machine learning model to predict the number of each amenity based on the number of residents within 1 km to 50 km at increments of 5 km. We use different machine learning models (linear regression, support vector machine, gradient boosting method, *k*-nearest neighbors regression and power law regression) and pick the best model based on test  $R^2$ .

Although many urban indicators do scale with population according to a power law relation<sup>29</sup>, by using machine learning we are able to exploit both linear and nonlinear relationships between population and amenities. This approach is further supported by having power law regression as the selected best model only for 17 (49%) out of 35 amenities.

## Results

### Normalization

Figure 1 clearly shows that percentile normalization is the best normalization method among the three methods that we have considered. Percentile normalization implies that the actual amenity count of a locality is not important—only their rank, which is quite understandable. When picking a locality among a set of candidates based on an amenity, a migrant does not really care

about the exact number of that amenity, only on which locality has the most of that amenity. This behavior has been hinted before by Noulas et al.<sup>26</sup>.

Percentile normalization also handles well amenity counts since their distribution is usually heavy-tailed and heavily skewed. Being bounded from zero to one, the weights are readily interpretable and comparable across models.

Min-max normalization models also performed better than RM. It is also bounded from zero to one so the weights are readily interpretable and comparable across models. Although the values were first transformed by  $x \rightarrow \log(1+x)$ , the absolute counts still matter and may have caused the poorer performance compared to percentile.

The adjusted  $z$ -score models performed worst, even worse than RM. It is also more difficult to compare and interpret because it is only bounded to the left by 0 but is not unbounded to the right.

Due to the consistent superior performance of percentile normalization, succeeding results and discussion will focus on percentile normalization models. A more detailed investigation of features were performed only on percentile normalization models as well.

## Feature sets

Population density seems to be a better feature compared to population. Switching RM to use population density (No normalization, Pop. dens. only in Fig. 1) instead of population results in a 2.5% relative improvement in MAPE. Looking at Fig. 1, we see that models with population density consistently outperform those with population albeit by a very small amount (0.056% to 5.9% relative improvement in MAPE) in many cases.

Combining population, population density and amenity counts do not result in the best model either. However, there is a caveat that the MAPE of the best four models are very close to each other (Amenities only MAPE = 6.1929%; Pop., Pop. dens. w/ amenities MAPE = 6.2014%). The best performing model is the one that only uses amenity counts implying that population information causally related to amenities and hence by virtue of granularity captures not just the mean field but also the variability, thereby providing a better model for human mobility.

Each OSM feature belongs to a primary group ([https://wiki.openstreetmap.org/wiki/Map\\_features](https://wiki.openstreetmap.org/wiki/Map_features)) such as building, landuse and shop. Another approach for grouping features is by performing Ward's agglomerative clustering<sup>30</sup> of the features based on the locality values. With these feature groupings, we investigate three methods of reducing the amenity-related features: (1) aggregating amenity count by primary features (*primary amenities*), (2) aggregating amenity count by groups based on feature clustering (*merged amenities*) and (3) selecting a representative amenity for each group based on feature clustering (*filtered amenities*). Reducing the number of amenity-related features results in worse performance compared to the best model but is still better than RM. Among the three methods considered, selecting a representative amenity for each group has the best performance (MAPE = 6.2%) whilst aggregating amenity counts by group performed worst (MAPE = 6.4%).

## Feature importance

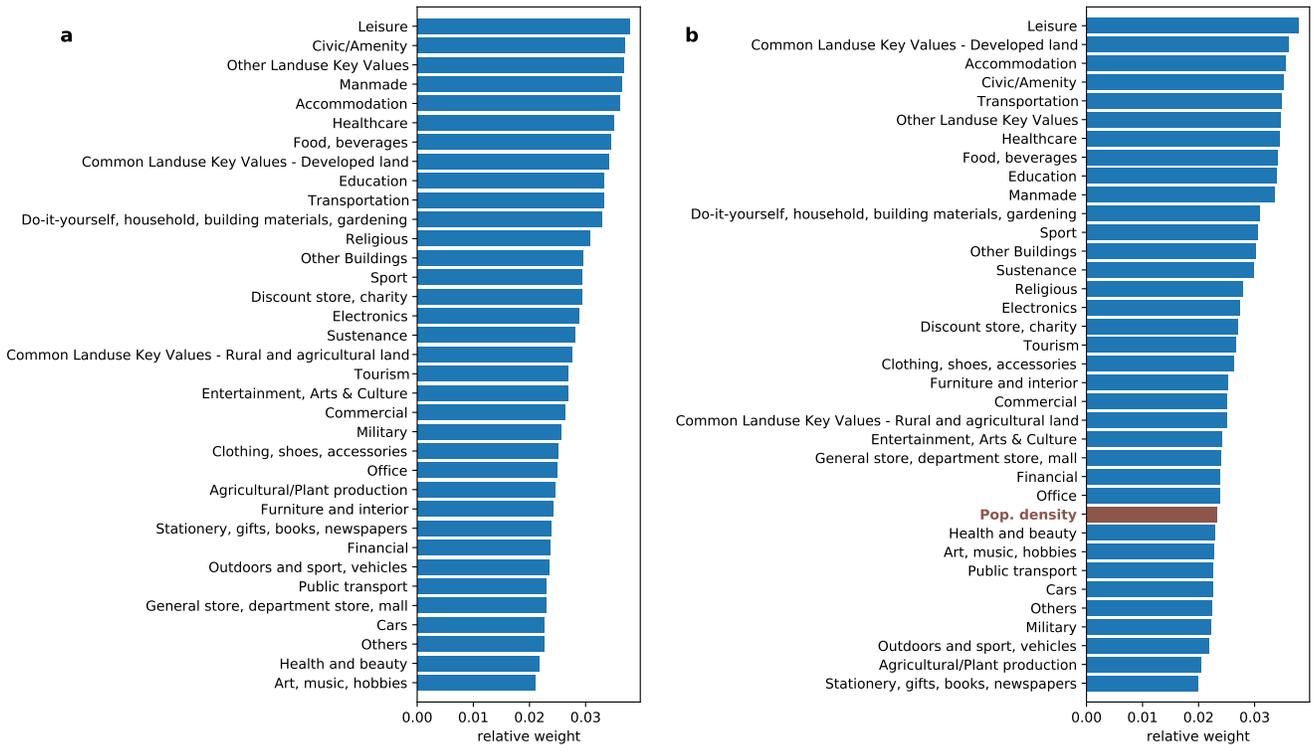
The relative weights of the amenity features of the best performing model (Percentile normalization, amenities only) are shown in Fig. 2a. At the outset, the top features do not seem to be related to jobs or work, supporting the initial assertion that migrants may not be solely looking at job opportunities when deciding to move; they may look at other concerns e.g., quality of life as well. However, this observation needs further investigation since more leisure and amenity places correspond to more service sector jobs. The presence of the top features may also correspond to the urbanity of a locality which could, in turn, imply more jobs.

The relative weights of features of the second best performing model (Percentile normalization, population density with amenities) still put Leisure at the top (Fig. 2b), however, there are several changes in the ranking of the features. This implies that small differences in the relative weight do not matter. Population density is the 27th out of 36 features in terms of relative weight, which suggests that most amenity features have higher relative weight than population information.

## Discussion

Predicting from and to where people move, and by how much, is one of the fundamental problems in the study of human mobility. RM provides a useful model for human mobility that allows us to answer this fundamental problem.

We extend this model by allowing amenities to be proxies for the migration attractiveness of a locality instead of population alone as in the original model. The model complements an earlier work demonstrating that amenities predict accurately the daily movement of people<sup>31</sup>. The result of the formulation shown here is consistent with how daily unchanging routines eventually accumulates to years resulting in permanent migration in some portion of the population. The model carries with it a natural way of interpreting the driver of migration to the level of amenities not possible in RM. Moreover it allows actionable insights that take into account the sensibilities or cultural preferences of the citizens of a country.



**Figure 2.** Feature importance of the two best performing models. (a) Percentile normalization, amenities only (b) Percentile normalization, population density with amenities. The most important features are not directly related to job opportunities which suggests people move not just because of job opportunities. Population feature importance is not ranked highly, even omitted in the best performing model which suggests amenities already include information derived from population.

Our results show that our model outperforms RM outright with as much as 7.7% relative improvement in MAPE for the best performing model. More importantly, amenity features outrank population features in importance with the best performing model not using any population feature at all. This suggests amenities already include information derived from population.

Amenities outweighing population information in terms of feature importance offers a couple of tantalizing applications. The first application is that this can potentially be used for doing population counts (census) in an area. With the feature weights as a guide, we can potentially investigate causality of amenities i.e., by how much will people be attracted to move to a place if we put up a particular amenity there. Of course, by doing so, we would be able to answer the conundrum of whether putting a particular amenity drives people to move there or is it the other way around—an amenity is put up because there are people there.

Generalizing and allowing a better resolved RM is a step closer in understanding more accurately the science of the emergence of cities. While the organization of amenities have been previously presented based on opposing concepts of diffusion and aggregation<sup>32,33</sup>, the complexity of the drivers that balance the built up and growth of cities are still an open concern<sup>34</sup>. The work here provides a procedure of quantifying a critical component of the formation of cities which is the movement of people as a function of diversity and quantity of amenities.

Our work is extensive: we considered different methods of normalization, feature sets, optimization target, feature weight estimator and even performance metric—distilling the results to only elaborate on the better performing configurations in this paper. It also provides an example of how machine learning can help resolve seemingly circular problems. In particular, using population to estimate amenity counts which will then be used for predicting population seems circular. However, by using machine learning, we were able to break this circular problem by incorporating nonlinearities that are not included in the power law model, which is the best theoretical model for the relationship between population and amenities. This approach also resulted in improved prediction, beating power law model 51% of the time.

## Data availability

All source data are publicly available at OSM (openstreetmap.org), GADM (gadm.org) and Philippine Statistical Authority (psa.gov.ph).

## References

1. Zhang, X., Xu, Y., Tu, W. & Ratti, C. Do different datasets tell the same story about urban mobility — A comparative study of public transit and taxi usage. *J. Transp. Geogr.* **70**, 78–90, DOI: [10.1016/j.jtrangeo.2018.05.002](https://doi.org/10.1016/j.jtrangeo.2018.05.002) (2018).
2. Piovani, D., Arcaute, E., Uchoa, G., Wilson, A. & Batty, M. Measuring accessibility using gravity and radiation models. *Royal Soc. Open Sci.* **5**, 171668, DOI: [10.1098/rsos.171668](https://doi.org/10.1098/rsos.171668) (2018). Publisher: Royal Society.
3. Tizzoni, M. *et al.* On the Use of Human Mobility Proxies for Modeling Epidemics. *PLOS Comput. Biol.* **10**, e1003716, DOI: [10.1371/journal.pcbi.1003716](https://doi.org/10.1371/journal.pcbi.1003716) (2014). Publisher: Public Library of Science.
4. Bengtsson, L. *et al.* Using Mobile Phone Data to Predict the Spatial Spread of Cholera. *Sci. Reports* **5**, 8923, DOI: [10.1038/srep08923](https://doi.org/10.1038/srep08923) (2015). Number: 1 Publisher: Nature Publishing Group.
5. Wesolowski, A., O'Meara, W. P., Eagle, N., Tatem, A. J. & Buckee, C. O. Evaluating Spatial Interaction Models for Regional Mobility in Sub-Saharan Africa. *PLOS Comput. Biol.* **11**, e1004267, DOI: [10.1371/journal.pcbi.1004267](https://doi.org/10.1371/journal.pcbi.1004267) (2015). Publisher: Public Library of Science.
6. Marshall, J. M. *et al.* Mathematical models of human mobility of relevance to malaria transmission in Africa. *Sci. Reports* **8**, 7713, DOI: [10.1038/s41598-018-26023-1](https://doi.org/10.1038/s41598-018-26023-1) (2018). Number: 1 Publisher: Nature Publishing Group.
7. Bagrow, J. P., Wang, D. & Barabási, A.-L. Collective Response of Human Populations to Large-Scale Emergencies. *PLOS ONE* **6**, e17680, DOI: [10.1371/journal.pone.0017680](https://doi.org/10.1371/journal.pone.0017680) (2011). Publisher: Public Library of Science.
8. Rutherford, A. *et al.* Limits of social mobilization. *Proc. Natl. Acad. Sci.* **110**, 6281–6286, DOI: [10.1073/pnas.1216338110](https://doi.org/10.1073/pnas.1216338110) (2013). Publisher: National Academy of Sciences Section: Physical Sciences.
9. Barbosa, H. *et al.* Human mobility: Models and applications. *Phys. Reports* **734**, 1–74, DOI: [10.1016/j.physrep.2018.01.001](https://doi.org/10.1016/j.physrep.2018.01.001) (2018).
10. Simini, F., González, M. C., Maritan, A. & Barabási, A.-L. A universal model for mobility and migration patterns. *Nature* **484**, 96–100, DOI: [10.1038/nature10856](https://doi.org/10.1038/nature10856) (2012). Number: 7392 Publisher: Nature Publishing Group.
11. Masucci, A. P., Serras, J., Johansson, A. & Batty, M. Gravity versus radiation models: On the importance of scale and heterogeneity in commuting flows. *Phys. Rev. E* **88**, 022812, DOI: [10.1103/PhysRevE.88.022812](https://doi.org/10.1103/PhysRevE.88.022812) (2013). Publisher: American Physical Society.

12. Yang, Y., Herrera, C., Eagle, N. & González, M. C. Limits of Predictability in Commuting Flows in the Absence of Data for Calibration. *Sci. Reports* **4**, 1–9, DOI: [10.1038/srep05662](https://doi.org/10.1038/srep05662) (2014). Number: 1 Publisher: Nature Publishing Group.
13. Kang, C., Liu, Y., Guo, D. & Qin, K. A Generalized Radiation Model for Human Mobility: Spatial Scale, Searching Direction and Trip Constraint. *PLOS ONE* **10**, e0143500, DOI: [10.1371/journal.pone.0143500](https://doi.org/10.1371/journal.pone.0143500) (2015). Publisher: Public Library of Science.
14. Amini, A., Kung, K., Kang, C., Sobolevsky, S. & Ratti, C. The impact of social segregation on human mobility in developing and industrialized regions. *EPJ Data Sci.* **3**, 1–20, DOI: [10.1140/epjds31](https://doi.org/10.1140/epjds31) (2014). Number: 1 Publisher: SpringerOpen.
15. Tigno, J. Migration and violent conflict in Mindanao. *Popul. Rev.* **45**, DOI: [10.1353/prv.2006.0013](https://doi.org/10.1353/prv.2006.0013) (2006).
16. Sterkens, C., Camacho, A. Z. & Scheepers, P. Ethno-religious Identification and Latent Conflict: Support of Violence among Muslim and Christian Filipino Children and Youth. In Harker, C., Hörschelmann, K. & Skelton, T. (eds.) *Conflict, Violence and Peace*, 1–16, DOI: [10.1007/978-981-4585-98-9\\_12-1](https://doi.org/10.1007/978-981-4585-98-9_12-1) (Springer Singapore, Singapore, 2016).
17. Chen, Y. & Rosenthal, S. S. Local amenities and life-cycle migration: Do people move for jobs or fun? *J. Urban Econ.* **64**, 519–537, DOI: [10.1016/j.jue.2008.05.005](https://doi.org/10.1016/j.jue.2008.05.005) (2008).
18. De la Roca, J. Selection in initial and return migration: Evidence from moves across Spanish cities. *J. Urban Econ.* **100**, 33–53, DOI: [10.1016/j.jue.2017.04.004](https://doi.org/10.1016/j.jue.2017.04.004) (2017).
19. Brown, D. L. & Wardwell, J. M. (eds.) *New Directions in Urban–Rural Migration: The Population Turnaround in Rural America* (Academic Press, 1980). Google-Books-ID: 63WLBQAAQBAJ.
20. Milbourne, P. Re-populating rural studies: Migrations, movements and mobilities. *J. Rural. Stud.* **23**, 381–386, DOI: [10.1016/j.jrurstud.2007.04.002](https://doi.org/10.1016/j.jrurstud.2007.04.002) (2007).
21. Stockdale, A. Contemporary and ‘Messy’ Rural In-migration Processes: Comparing Counterurban and Lateral Rural Migration. *Population, Space Place* **22**, 599–616, DOI: <https://doi.org/10.1002/psp.1947> (2016). \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/psp.1947>.
22. Liu, E. & Yan, X. New parameter-free mobility model: Opportunity priority selection model. *Phys. A: Stat. Mech. its Appl.* **526**, 121023, DOI: [10.1016/j.physa.2019.04.259](https://doi.org/10.1016/j.physa.2019.04.259) (2019).
23. Liu, E.-J. & Yan, X.-Y. A universal opportunity model for human mobility. *Sci. Reports* **10**, 4657, DOI: [10.1038/s41598-020-61613-y](https://doi.org/10.1038/s41598-020-61613-y) (2020). Number: 1 Publisher: Nature Publishing Group.
24. Robinson, C. & Dilkina, B. A Machine Learning Approach to Modeling Human Migration. In *Proceedings of the 1st ACM SIGCAS Conference on Computing and Sustainable Societies, COMPASS '18*, 1–8, DOI: [10.1145/3209811.3209868](https://doi.org/10.1145/3209811.3209868) (Association for Computing Machinery, New York, NY, USA, 2018).
25. McCulloch, K., Golding, N., McVernon, J., Goodwin, S. & Tomko, M. Ensemble model for estimating continental-scale patterns of human movement: a case study of Australia. *Sci. Reports* **11**, 4806, DOI: [10.1038/s41598-021-84198-6](https://doi.org/10.1038/s41598-021-84198-6) (2021). Number: 1 Publisher: Nature Publishing Group.
26. Noulas, A., Scellato, S., Lambiotte, R., Pontil, M. & Mascolo, C. A Tale of Many Cities: Universal Patterns in Human Urban Mobility. *PLOS ONE* **7**, e37027, DOI: [10.1371/journal.pone.0037027](https://doi.org/10.1371/journal.pone.0037027) (2012). Publisher: Public Library of Science.
27. Stouffer, S. A. Intervening Opportunities: A Theory Relating Mobility and Distance. *Am. Sociol. Rev.* **5**, 845–867, DOI: [10.2307/2084520](https://doi.org/10.2307/2084520) (1940). Publisher: [American Sociological Association, Sage Publications, Inc.].
28. Barrington-Leigh, C. & Millard-Ball, A. The world’s user-generated road map is more than 80% complete. *PLOS ONE* **12**, e0180698, DOI: [10.1371/journal.pone.0180698](https://doi.org/10.1371/journal.pone.0180698) (2017). Publisher: Public Library of Science.
29. Bettencourt, L. M. A., Lobo, J., Helbing, D., Kühnert, C. & West, G. B. Growth, innovation, scaling, and the pace of life in cities. *Proc. Natl. Acad. Sci.* **104**, 7301–7306, DOI: [10.1073/pnas.0610172104](https://doi.org/10.1073/pnas.0610172104) (2007). Publisher: National Academy of Sciences Section: Social Sciences.
30. Aggarwal, C. C. *Data mining: the textbook* (Springer, 2015).
31. Hu, N., Legara, E. F., Lee, K. K., Hung, G. G. & Monterola, C. Impacts of land use and amenities on public transport use, urban planning and design. *Land Use Policy* **57**, 356–367, DOI: [10.1016/j.landusepol.2016.06.004](https://doi.org/10.1016/j.landusepol.2016.06.004) (2016).
32. Decraene, J., Monterola, C., Lee, G. K. K. & Hung, T. G. G. A Quantitative Procedure for the Spatial Characterization of Urban Land Use | International Journal of Modern Physics C. *Int. J. Mod. Phys. C* **24**, DOI: <https://doi.org/10.1142/S0129183112500921> (2013).

33. Decraene, J., Monterola, C., Lee, G. K. K., Hung, T. G. G. & Batty, M. The Emergence of Urban Land Use Patterns Driven by Dispersion and Aggregation Mechanisms. *PLOS ONE* **8**, 1–9, DOI: [10.1371/journal.pone.0080309](https://doi.org/10.1371/journal.pone.0080309) (2013). Publisher: Public Library of Science.
34. Ortman, S. G., Lobo, J. & Smith, M. E. Cities: Complexity, theory and history. *PLOS ONE* **15**, e0243621, DOI: [10.1371/journal.pone.0243621](https://doi.org/10.1371/journal.pone.0243621) (2020). Publisher: Public Library of Science.

## **Acknowledgements**

We thank Antonino Paguirigan, Jr. and Eduardo David, Jr. for initial discussions. This work is primarily funded by the DOST-PCIEERD with Project No. 08501, 2020. The funding source has no involvement in the conduct of the research and preparation of this article.

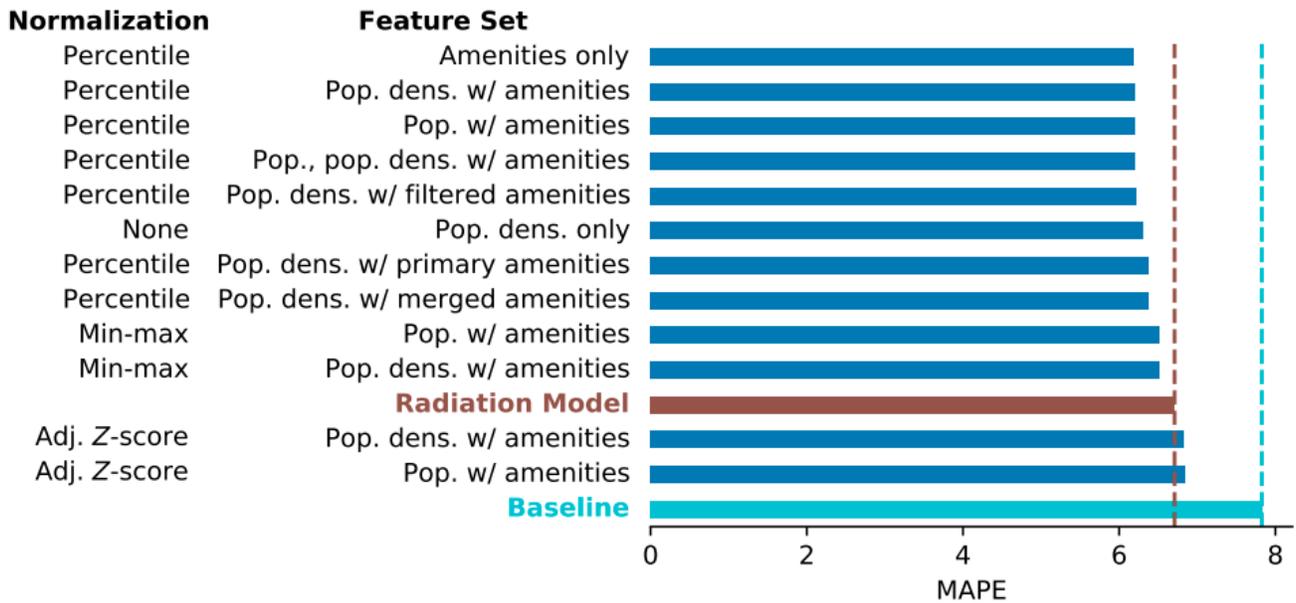
## **Author contributions statement**

All authors conceived and designed the study, analysed and interpreted the results. C.M.A. gathered and manipulated the data, carried out the analysis, wrote all the codes and the first draft. All authors reviewed the manuscript.

## **Competing interests**

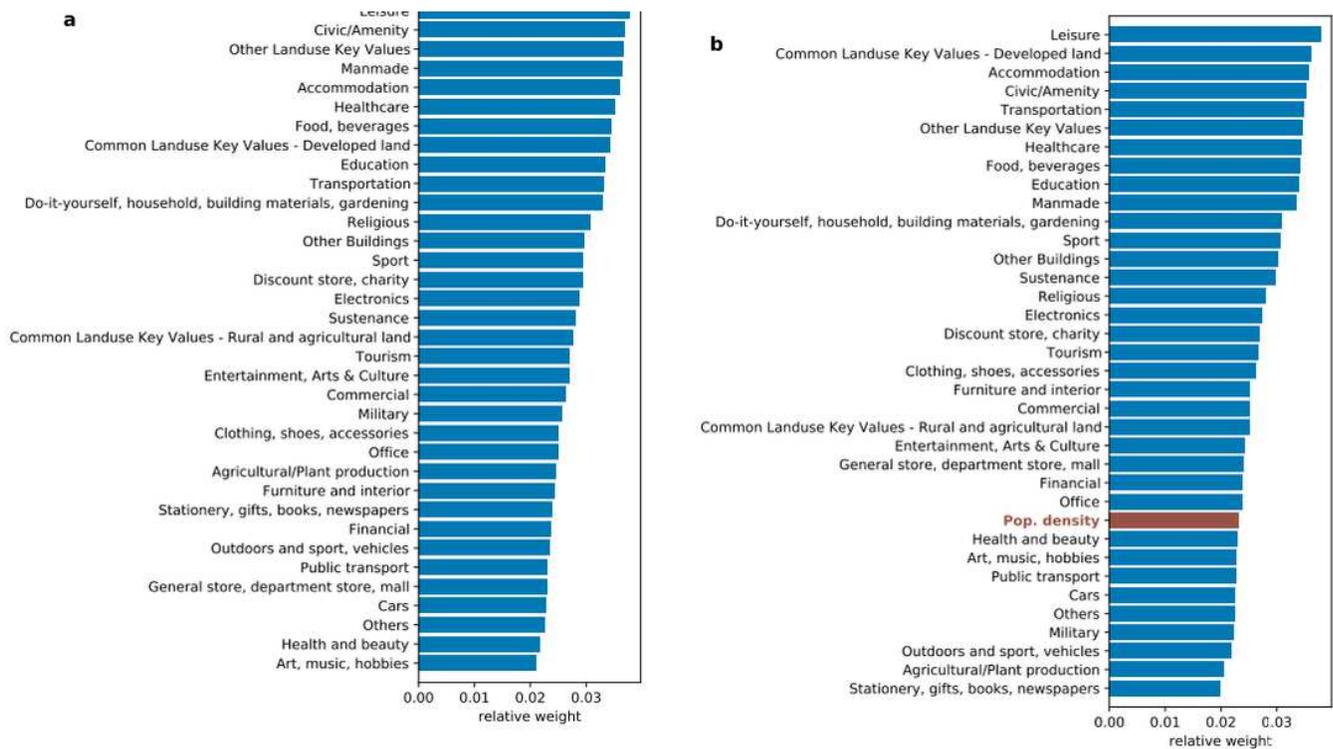
The authors declare no competing interests.

# Figures



**Figure 1**

Model performance based on mean absolute percentage error (MAPE). Percentile normalization is the best method of normalization, outperforming Min-max and Adjusted Z-score normalizations in all instances. The best performing model uses only amenities as features and follows Percentile normalization. It corresponds to a 7.7% MAPE improvement relative to Radiation Model. All models beat the baseline model which is the outright scaling of locality population according to the same rate of change in the national population. The performance metrics of all of the 78 configurations that were investigated are displayed in Supplementary Table S1.



**Figure 2**

Feature importance of the two best performing models. (a) Percentile normalization, amenities only (b) Percentile normalization, population density with amenities. The most important features are not directly related to job opportunities which suggests people move not just because of job opportunities. Population feature importance is not ranked highly, even omitted in the best performing model which suggests amenities already include information derived from population.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [si.pdf](#)