

Do road-traffic injuries scale non-linearly with travel?

Rob Johnson | <https://github.com/robj411> | rob.johnson@mrc-bsu.cam.ac.uk

November 12, 2019

Contents

1	Introduction	3
2	Background	5
2.1	Equation 2 in the literature	5
2.2	Illustration with data for England	5
2.3	The scope of this work	6
3	Examination of Equation 2	7
3.1	Interpretation of Equation 2	7
3.1.1	Methodology applied to infectious disease	7
3.2	Application of Equation 2	8
3.3	Scalability of Equation 2	9
3.3.1	Illustration: data observed over different periods of time	9
3.3.2	Example: England injury data	9
4	Use of learnt coefficients in predictive models for new settings	13
4.1	Scaling	13
4.2	Mode definitions	15
4.3	Number of modes	16
5	Simulation and motivation for an alternative model	17
5.1	Simulation model	17
5.2	Simulation study	17
5.3	Comparison of model predictions to simulations	18
5.4	Simulated safety in numbers	19
5.5	Discussion	21
5.6	Conclusion	21
6	An alternative model	23
6.1	The model	23
6.2	Discussion	24
7	A size-adjusted study of the England data	25
7.1	Results	25
7.2	Discussion	26
8	Conclusions	28

A	Supplementary figures	30
B	Worked example: implementation in ITHIM-R	32
B.1	Constructing the model	32
B.2	Making predictions	32
B.3	β_1 and β_2 parameters	32
C	Power-law scaling relationships	34

Abstract

We consider the mathematical model used to describe the phenomenon of “safety in numbers” in relation to road-traffic injuries. We focus on a model commonly used to describe number of injuries as a function of total use of two transport modes – motor cars and cyclists – in large contiguous geographic areas, such as cities and counties. While models often include many covariates, the basic functional form does not account for the size of the area under study (or, equivalently, the density of road use per unit time and/or space). We present some consequences of this modelling choice in terms of interpretation of meaning and application in predictive modelling, and explore it further through simulations and analysis. We develop the simulation and analytic results into a new model that we demonstrate using data from England. We conclude that omission of the size of an area leads to confounding which results in a bias in the model, rendering it inadequate to make predictions of change in injury number following a change in transport mode use.

1 Introduction

The term “safety in numbers” reflects the observation that a change in the number of road users of a specific type (commonly cyclists) is not met by a proportional change in the number of injuries caused to or by the road-user group. Equivalently, from an individual-level perspective, the more cyclists there are, the lower the risk is per cyclist. We describe this as “linearity”: is the expected (total) number of injuries a linear function of the number of cyclists (with an intercept at zero)? If it is “sublinear”, we might say that there is safety in numbers. Our challenge is to formulate this question in a mathematical model and take account of all relevant covariates.

The number of injuries over a certain period and space, I , is generally formulated as follows, in terms of a base rate, α , the number of road users of one type (say, motorists), M , the number of road users of the second type (say, cyclists), C , and “safety in numbers” exponents for each, β_1 and β_2 (Elvik and Bjørnskau, 2017):

$$I \sim \mathcal{F}(\lambda), \quad (1)$$

$$\lambda = \alpha M^{\beta_1} C^{\beta_2}, \quad (2)$$

where \mathcal{F} is a function specifying a distribution, such as Poisson or negative binomial. For completeness, one might include other covariates, e.g.

$$\lambda = \alpha M^{\beta_1} C^{\beta_2} \exp \left(\sum_{i=3}^P \beta_i X_i \right), \quad (3)$$

but these are not central to the present discussion, which focuses around the variables in Equation 2.

Beyond descriptions of observed data, such analyses are proposed to inform public-health forecasting (Schepers and Heinen, 2013) and to make predictions in novel scenarios (de Sá et al., 2017). This includes assessment of likely health benefits following policy change, as well as forecasting healthcare needs given the expected change to transport-related behaviour, e.g. the increase in motor vehicle ownership expected in cities in fast-growing economies in the coming years, especially in light of Sustainable Development Goal 3.6 to reduce road-traffic casualties to half of their 2011 value.¹ We have previously used this formulation to make predictions of numbers of road injuries for specific cities (Accra, Sao Paulo), as well as in the developing generic software ITHIM-R. (ITHIM: integrated transport and health impact model.) To make predictions in novel scenarios requires an assumption of causality, and the validity of a causal relationship relies on having accounted for all confounders. [This, in turn, relies on assumptions of “positivity” (existence of sufficient data informing the contrast with/without the intervention) and “exchangeability” (confounding accounted for) of the data.]

In this work, we question the validity of the assumptions implicit in the model as well as its applications. Our intention here is to examine (1) the interpretation of (and language around)

¹<http://iris.wpro.who.int/handle/10665.1/12878>

Equation 2, (2) its proposed application to health-impact modelling, and (3) what “safety in numbers” is and how it might be measured. We might consider two types of study: studies comparing multiple roads or road junctions in the same city, and studies comparing different cities (or geographical regions containing whole road networks). In this work we focus on the latter, but make reference to the former, as it provides some insights, it might be used for prediction, and in the end we would like to have a single comprehensive framework. The picture of questions is shown in Table 1.

Table 1: Questions & answers for inter-city studies

Question	Answer	See Section
(1) What does it mean?	$\beta_1 + \beta_2 = 1$ means linearity	3.1
(2) How do we use it to predict?	Account for density	5 and 7
(3) How do we study it?	Account for density	7 and Appendix C

Using a theoretical approach, we derive the result that coefficients $\beta_1 + \beta_2 = 1$ in Equation 2 learnt from multiple settings, which differ in scale but have the same road-user density, correspond to linear scaling across time and space. Put another way, these coefficients permit “tiling” of a small space to create a larger space whose properties are the same in terms of risk per unit space. We confirm this result through simulation, which shows also that, again under Equation 2, the coefficients $\beta_1 + \beta_2 = 2$ will tend to describe the variations seen in injury count data from multiple areas of the same size but different road-user density. These provide null hypotheses that injuries are linear in road use for studies assessing the impact of road-user number on collision risk. It is useful to understand “tiling” in contrast with mode shifts: tiling preserves properties across space and time, allowing us to make predictions for longer or shorter time periods, or larger or smaller areas, where we expect road use to follow scale. We contrast this with predictions where the time and space stays the same and the road use changes.

We present an alternative model which includes the components of Equation 2 and additionally it explicitly includes size, since it is a cause of I , M and C , and is therefore a confounding factor for estimating, from data from areas of different sizes, the effect of interventions that change M and C . (Note that size might be parametrised in terms of area, population, road length, or any other relevant metric.) We use exponents δ_1 and δ_2 to parametrise this model, to distinguish it from existing models. The cases discussed above (scaling size, constant density ($\beta_1 = \beta_2 = 0.5$) and constant size, scaling density ($\beta_1 = \beta_2 = 1$)) are special cases of this model. We apply this model to data for England. Our results are consistent with the preceding theoretical and simulation analyses. With this model and these data, we reject the null hypothesis of linearity of injuries with respect to road-user density for most subsets of the data – that is, for doubled travel, we expect fewer than twice the number of injuries.

The main purposes of this work are (a) to identify an acceptable model for prediction of injury numbers under a shift in road use, and (b) to open new avenues for research into road-traffic injury dynamics. We highlight some important features of existing the existing model and propose an amended model for prediction. We identify areas we believe are most in need of attention, which are: accounting for size/density; the definition of “density”; the link between small-scale and city-level studies; how we interpret non-linearity where a predictor is a group of modes; and how we identify “non-linearity”.

2 Background

2.1 Equation 2 in the literature

“Safety in numbers” exponents are estimated in analyses of road-traffic injuries through fitting a regression model such as Equation 2 or Equation 3 to data. At minimum, these data consist of counts of road injuries, and of two road user types, often cars and pedestrians or cyclists. The exponents β_1 and β_2 are estimated, and reported with confidence intervals and p-values corresponding to the probability of observing the data under the null hypothesis $\beta_1 = 0$, $\beta_2 = 0$, which is the default in most statistical software.

Various different units of measurement have been used in studies reporting safety-in-numbers effects. For example, Miranda-Moreno et al. (2011); Geyer et al. (2006); Garder et al. (1998); Schepers et al. (2011); Nordback et al. (2014) and Leden (2002) count vehicles (usually in terms of average daily number of vehicles per unit space, though sometimes the time unit is annual or hourly). In contrast, Prato et al. (2016) and Schepers and Heinen (2013) work in terms of km travelled by each road-user group. Injuries are counted as the total in one or more years. The areas considered in the studies range in size from intersections (Nordback et al., 2014) to municipalities (Schepers and Heinen, 2013) and “local authorities” (Aldred et al., 2017).

From this we see that time and space are not typically explicitly included in published models. Instead, it is implicit that the numbers of injuries are linear in time and space. We return to this issue later, where we see that this specification results in a base rate (α in Equation 2) that depends on e.g. the number of years of study (Section 3.3.1).

In terms of prediction, exponents can be used as “laws”, as in de Sá et al. (2017), where learnt exponents are applied alone, outside the context of the training data and the full model. This differs from what we usually understand of making predictions from models, in which all parameters from a model are used to compute its prediction(s). For ITHIM-R/Tigthat, as in de Sá et al. (2017), we take our current observations (obs) and make a prediction (pred) as follows:

$$I_{\text{pred}} = I_{\text{obs}} \frac{\alpha M_{\text{pred}}^{\beta_1} C_{\text{pred}}^{\beta_2}}{\alpha M_{\text{obs}}^{\beta_1} C_{\text{obs}}^{\beta_2}} = I_{\text{obs}} \left(\frac{M_{\text{pred}}}{M_{\text{obs}}} \right)^{\beta_1} \left(\frac{C_{\text{pred}}}{C_{\text{obs}}} \right)^{\beta_2}. \quad (4)$$

Note that the dependence on the base rate α is lost and injuries are predicted as a function of the old observed injury counts and travel, the new/predicted travel, and the parameters β . Similarly, were we to use a model of the form of Equation 3, the exponential terms would cancel out.

In this way, the method relates to power-scaling laws in the biological sciences, dating back to metabolic rate as a function of size (Kleiber, 1947), and extending to social and economic sciences (Bettencourt et al., 2007). Laws hypothesised resulting from models such as these have been widely discussed in terms of both methods and implications (Stumpf and Porter, 2012; Clauset et al., 2009), which have some relevance for our example.

2.2 Illustration with data for England

To give some context, we present data pertaining to road-traffic injuries recorded in England in the years 2005 to 2015 for 148 areas, including counties and boroughs. From among these data we isolate all injuries to cyclists that occurred in events involving at least one car. Alongside these data, we use Road Traffic Statistics estimates of distance travelled by cars and bikes in these areas², as well as census estimates for population numbers in the year 2011.

We begin with some descriptive figures demonstrating what the data look like. We show how injuries to cyclists (involving cars) scales with the total distance travelled by cyclists (Figure 1), demonstrating how “safety in numbers” might be identified through comparison of the gradient of the line of best fit to one. In this Figure, we show side by side how the gradient changes as we include only more severe cases. We repeat the analysis with London boroughs only in Figure 15.

²We use data covering all road types provided by RTS minus that on motorways, calculated from <https://roadtraffic.dft.gov.uk/downloads>.

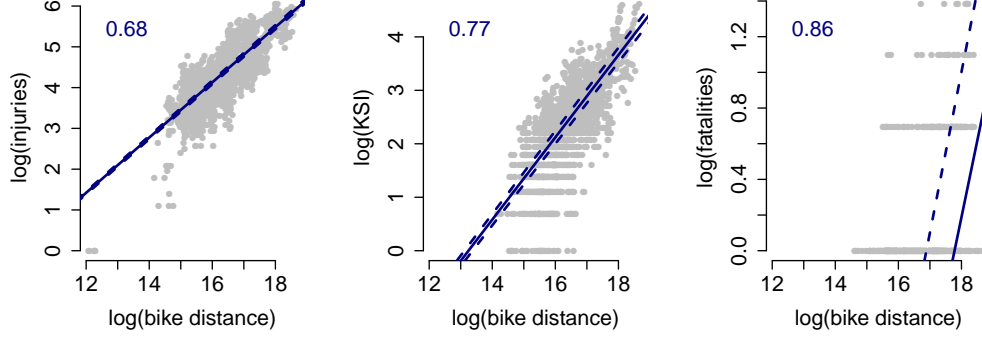


Figure 1: Fit of equation $I = \alpha C^\beta$ to data for local areas in England. I is (left) total incidents involving a car in which a cyclist was injured; (middle) the subset of these that were KSI; and (right) the subset of these that were fatal. C counts the total cycle distance travelled in the area. In blue is the line of best fit, with the gradient of the line (the value for β) in the top-left corner.

We look also at many possible relationships between single predictors and injury counts (Figure 16). Note the similar trend as the less serious casualties are excluded (KSI: killed or seriously injured). These values are summarised in Table 2, where we include also the result of a regression of the form of Equation 2. Note that the sum aligns with the other rows.

Table 2: Coefficients

Predictor	All injuries	KSI	Fatalities
Total travel (β)	0.77	0.90	1.11
Population (β)	0.78	0.89	0.98
Bike (β)	0.68	0.77	0.86
Car (β)	0.75	0.88	1.11
Bike+Car ($\beta_1 + \beta_2$)	0.67	0.79	0.97

2.3 The scope of this work

We focus our attention on city-level models (i.e. large contiguous geographic areas such as cities, counties and boroughs), rather than small-scale models. Our reasons for doing so are (a) brevity, (b) relevance to our aims (question 2 of Table 1), and (c) the ready availability of test data (Section 2.2). We would ultimately like to join up the work presented here with studies of small-scale areas, but more groundwork is needed first to understand how the features we explore relate to that setting, and what features we have omitted are of greater importance at a small scale.

Studies of units the size of cities are easier to conceptualise and, possibly, easier to model in an intuitive way. Relationships between size, travel, population, road length and other city-level metrics will be close to linear. Local features that might be significant on a smaller scale are likely averaged out at the level of a city, for example the distribution of travel over time.

Therefore, although we consider how small-scale studies might be used to make predictions for the city level (Section 4.1), we do not directly address questions 1 and 3 of Table 1 with reference to small-scale studies as further work in this area is required.

3 Examination of Equation 2

In general, it is posited that $\beta_i < 1$ for an equation of the form of Equation 2 implies safety in numbers for mode i (Elvik and Bjørnskau, 2017). However, given the construction of Equation 2, we show that $\beta_i < 1$ for all i for any study counting motorists, cyclists and injuries across scales, because the mode counts C and M are functions of scale. This has implications both for our interpretation and for our application of such a model.

3.1 Interpretation of Equation 2

To illustrate, consider again the example of the English counties, excluding London boroughs, and cyclist KSI counts. Applying Equation 2, we find coefficients $\beta_1 = 0.17$ and $\beta_2 = 0.63$, and we might interpret this as safety in numbers.

However, there is an approximately linear relationship between population (N) and cyclist KSI, between population and car distance, and between population and cyclist distance (Figure 17). If we consider population to be a confounder that explains the relationships between the other three variables, and fit the model $I = \alpha N^\beta$, we find $\beta = 0.96$ with standard error 0.03 and p value 0.28 – i.e., injuries are close to linear in population and, from these data, we can't rule out that they are linear.

The act of decomposing the population N into two component parts (M and C , where we have seen that there is a linear relationship between the three) splits the exponent associated with N between M and C . We would observe a three-way split were we to consider three modes.

3.1.1 Methodology applied to infectious disease

We would see the same phenomenon in HIV infection rates, if we were to count the number of new infections I and divide the population N into the susceptible group, C , and the infectious group, M (Bettencourt et al., 2007). We can examine this example to see what the implications of the regression method are on how we understand safety scaling as a function of population number.

Let's start with hypothetical, noiseless "data" that illustrate the relationship $I = \alpha C^{\beta_1} M^{\beta_2}$, yielding $\beta_1 = \beta_2 = 0.5$, exemplary safety in numbers (Table 3).

Table 3: A simple example of idealised data representing safety in numbers.

Cyclist injuries I	Motorists M	Cyclists C
1	100	10
2	200	20
3	300	30
4	400	40

Now let's relabel the columns and consider instead the number of new infections in a year (I), the number of infectious people (C) and the number of susceptible people (M), Table 4.

Table 4: A simple example of idealised data representing safety in numbers applied to infectious disease, where we count the number of new infections per year as a function of the number of susceptible and the number of infectious people.

New infections I	Susceptible M	Infectious C
1	100	10
2	200	20
3	300	30
4	400	40

So we see that for infectious diseases we also have safety in numbers. In fact, the rate of a disease (per capita) is often studied as a function of the total population, N , rather than the infectious and

susceptible populations. We see from Table 5 that infections are linear in population ($I = \alpha'N$), even though there is safety in numbers for both infectious and susceptible parties.

Table 5: A simple example of idealised data representing safety in numbers applied to infectious disease, where we count the number of new infections per year as a function of total population.

New infections I	Susceptible M	Infectious C	Population N
1	100	10	110
2	200	20	220
3	300	30	330
4	400	40	440

For HIV infections, a relationship of $I \sim \alpha'N^{1.2}$ has been observed (Bettencourt et al., 2007). Given that $I = \alpha'N$ corresponds to $I = \alpha M^{0.5}C^{0.5}$, as in Table 5, how would we write $I = \alpha M^{\beta_1}C^{\beta_2}$ for $I \sim \alpha'N^{1.2}$? Perhaps $I = \alpha M^{0.6}C^{0.6}$? Although this is pure speculation, it seems reasonable as a first guess. In any event, we would not expect either β_1 or β_2 to exceed 1. The resulting inference in the safety-in-numbers framework is then that the biggest susceptible populations M are the safest. However, the biggest susceptible populations are in the locations with the biggest populations, which have, according to the observation $I \sim \alpha'N^{1.2}$, the highest rates of new infections per year per capita.

3.2 Application of Equation 2

We ought also to consider the consequences of collinearity in using our models to make predictions. It is well known that the validity of predictive models worsens as one departs from the training space; all the more so with colinear variables, whose predictive performance is poor even within the training space (Kiers and Smilde, 2007). In addition, in many scenarios we are likely to consider in health-impact modelling, we are particularly interested in mode shifts (Schepers and Heinen, 2013): that is, we are considering transitions that run perpendicular to the training data used in the construction of the model, increasing one mode and decreasing another, rather than increasing or decreasing both together (Figure 18). We return to these dynamics with a particular focus on density in a simulation study in Section 5.

In summary, we would do well to express as much as we can of the known data-generating mechanism in the regression model, e.g. accounting for size properly. Then we would need not rely on questionable assumptions such as that there is no confounding.

3.3 Scalability of Equation 2

Equation 2 fit at one scale cannot be used to make a prediction at another scale. A model fit to parts of a whole cannot inform an unbiased prediction for the whole. The contradiction is in part due to working in numbers rather than density: the equation does not distinguish between differences in number due to extended measurement, i.e. recording injuries over a larger area or time period, and an increase in number due to there being more road users in the same area and time period, i.e. a greater density of road users. The consequence is that each parametrised model is particular to its own setting, which vary in scale from hours to years, and junctions to counties. Despite these ranges in scale, remarkably consistent values are reported for effects of “safety in numbers”. An illustration and an example are given in Sections 3.3.1 and 3.3.2.

In fact, the absence of accounting for scale might be what gives rise to the “safety in numbers” observation in the first place. Formally, we start with a single observation of injury number $I = \alpha M^{\beta_1} C^{\beta_2}$, and make a second observation identical to the first, which, added to the first, gives $2I = \alpha(2M)^{\beta_1} (2C)^{\beta_2}$. We can solve these equations together to learn about the β values where we have assumed linearity:

$$\frac{2I}{I} = \frac{\alpha(2M)^{\beta_1} (2C)^{\beta_2}}{\alpha M^{\beta_1} C^{\beta_2}}; \quad (5)$$

$$2 = \frac{2^{\beta_1} M^{\beta_1} 2^{\beta_2} C^{\beta_2}}{M^{\beta_1} C^{\beta_2}}; \quad (6)$$

$$= \frac{2^{\beta_1} 2^{\beta_2}}{1}; \quad (7)$$

$$= 2^{\beta_1 + \beta_2} \quad (8)$$

So $\beta_1 + \beta_2 = 1$ when injuries are linear in observation sizes. Note that this result is independent of the supposed or true mechanism, and independent of the relationship between the covariates M and C . Therefore, as $\beta_1 + \beta_2 = 1$ – safety in numbers – is consistent with a simple linear relationship between injuries the sizes of observations we must ask whether the results that have been observed arose as a result of such a construct in the dataset under study.

3.3.1 Illustration: data observed over different periods of time

We simulate a study of observations made in a single area over different numbers of years, but where injury risk is the same over time. Starting with Equation 2, we set $\alpha = \exp(-6.879)$, $\beta_1 = 0.591$, and $\beta_2 = 0.32$, and simulate data for a single study where M takes mean value 8000 and C takes mean value 340. To simulate many studies that vary by duration, we simply multiply the data by the size of the study. We learn the parameters through a regression model corresponding to Equation 2.

Figure 2 shows that, as the size of the study increases, the base rate, α , also increases, even though the individual-level injury risks underlying the data are exactly the same. This identifies the trade-off between the parameters β_i and α : if β_1 and β_2 are fixed, our base rate α will increase as the size of the study increases. My intuition is that, in this trade-off, it is the rate α that should stay constant, while the exponents might vary with study size (Figure 3).

The purpose of this illustration is to show that if the data are generated by a mechanism such that the data points are distinguished by different observation sizes, but the risk is constant over all data points, then fitting Equation 2 is inappropriate for predicting risk. The model is mis-specified, confounded as it is by observation size, which is a cause of both M and C . Omission of size from the regression causes bias in the estimated baseline risk and the estimated effects of M and C , and is therefore inappropriate for inferring a universal effect and predicting effects of mode shift.

3.3.2 Example: England injury data

Our example uses the STATS19 injury data for England in the years 2005 to 2015, across 148 areas (which include counties and local authorities). Areas are grouped into nine regions. We fit Equation 2 (a) to area-level data within regions and make predictions for each region, given the

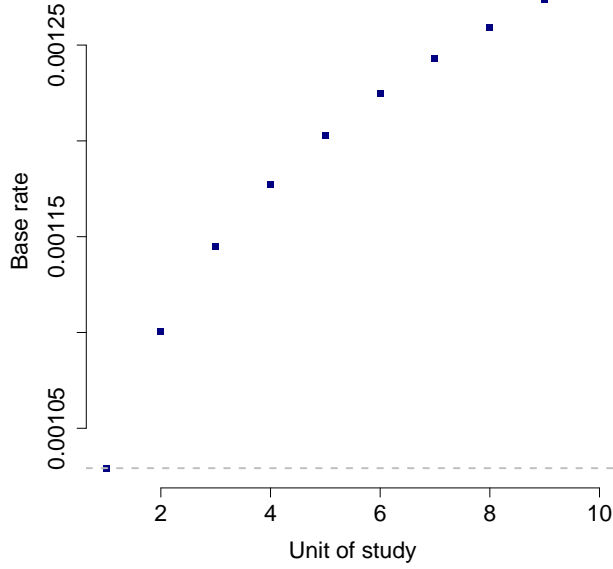


Figure 2: Base rate (α) as the size of a study increases, and safety-in-numbers exponents are fixed. The grey dashed line shows the true base rate.

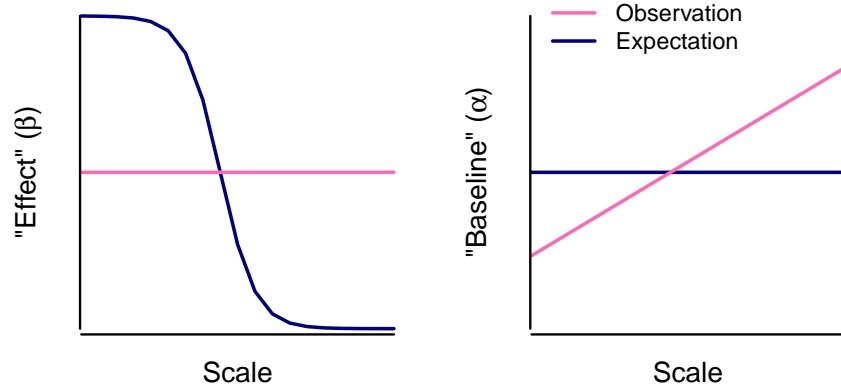


Figure 3: Safety-in-numbers effect (β) and base rate (α) as the size of a study increases. Studies observe the effect to stay more or less constant across scales, while the base rate increases (pink). One might expect instead the effect to diminish across scales and the base rate to stay constant (navy blue).

observed values of M and C in each region, and (b) to area-level data for the whole country and make predictions for the country.

Figure 4 shows the coefficients we infer for cyclist KSI counts resulting from collisions with other cyclists, motorcyclists, cars, vans, buses and heavy-goods vehicles. Empty circles are region estimates based on the areas they contain. In orange are the country estimates using all areas. In turquoise are the country estimates using the nine regions. Note the trend of clustering towards the line $\beta_1 + \beta_2 = 1$.

We now use the coefficients to make predictions at higher scales (Figure 5). There is a bias in underestimating mixed-mode injuries and overestimating same-mode injuries. In particular, we observe $\sim 16,000$ injuries to cyclists caused by cars in the country data, but predict only $\sim 8,000$.

Many studies have estimated similar values for β_1 and β_2 , and this is taken as evidence for safety

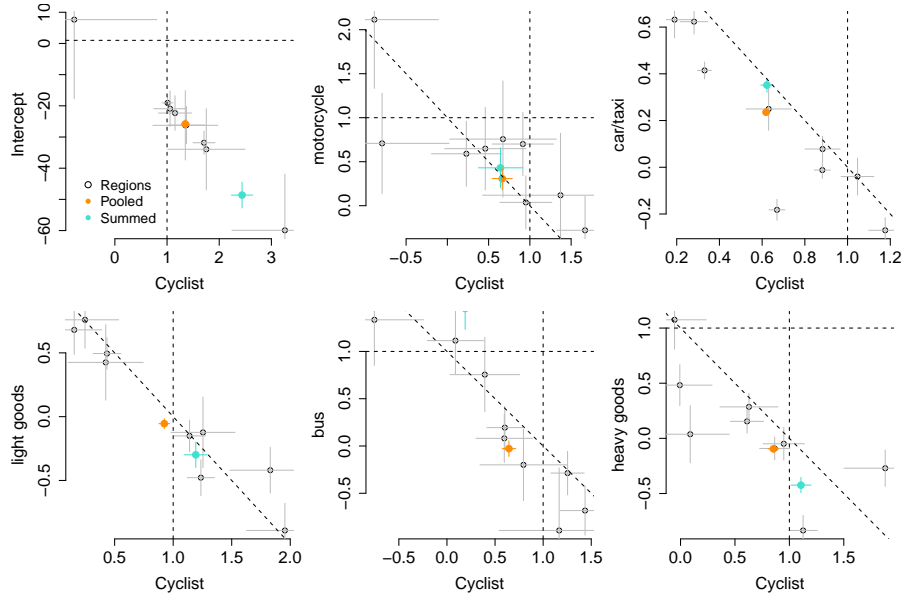


Figure 4: Coefficients we infer for cyclist KSI counts resulting from collisions with other cyclists, motorcyclists, cars, vans, buses and heavy-goods vehicles. Empty circles are region estimates based on the counties they contain. In orange are the country estimates using all counties. In turquoise are the country estimates using the nine regions. Note that for cyclist–cyclist collisions, there is only one coefficient, so we plot this coefficient against the intercept.

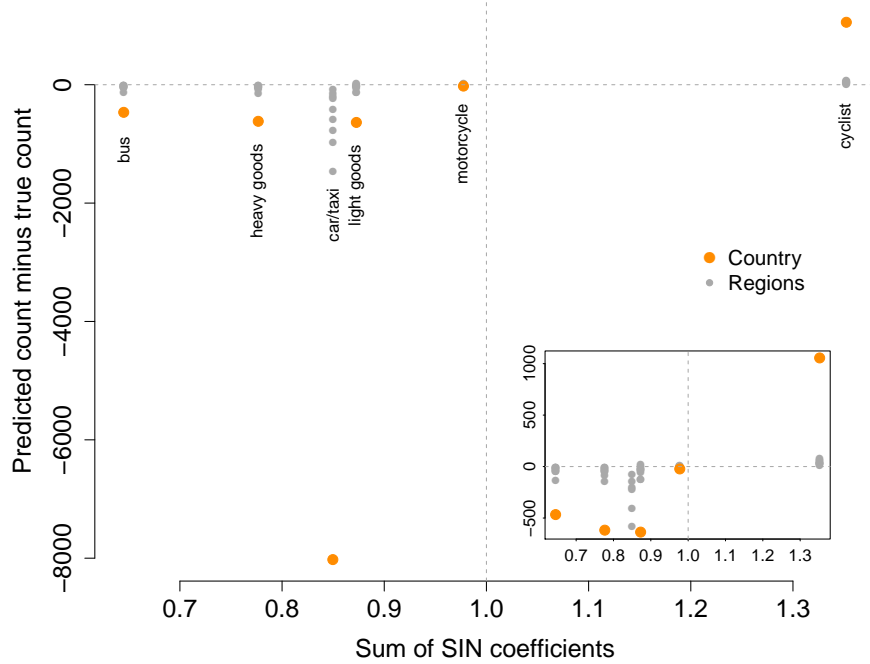


Figure 5: Predictions we make for cyclist KSI counts relative to observed counts, resulting from collisions with other cyclists, motorcyclists, cars, vans, buses and heavy-goods vehicles.

in numbers. Note, however, that studies examined areas with different scales, and the effect does not diminish as size increases, as we have shown in Section 3.3 (in fact it tends to $\beta_1 + \beta_2 = 1$). Therefore, the effect identified is not a function of proximity. If our observed values indicating

safety in numbers are not to do with locality, then where do they come from? Might it be that it is not an effect of safety at all, but rather an artefact of omission of scale from the model?

4 Use of learnt coefficients in predictive models for new settings

We are particularly interested in applying learning from regression modelling to making predictions under interventions or scenarios in which the numbers of motorists and cyclists change, for which we assume that the observed safety-in-numbers relationship is causal. First, we explore the possibility of applying coefficients learnt at smaller scales to city-level prediction models. We describe a general method and discuss issues remaining to be resolved.

4.1 Scaling

Ideally we would be able to collect data for the setting we are interested in, and build a model from scratch that includes all potential confounders, but in most circumstances that won't be possible. However, previous studies will contain “indirect” information, i.e. partially relevant settings. How can we use the results of studies such as those documented in Elvik and Bjørnskau (2017) in a meaningful, appropriate way in city-level prediction models? (N.B.: Of interest to us is the non-linearity of the relationship encapsulated in “safety-in-numbers coefficients” – the intercepts and other covariates we can learn independently for our setting.) We can define an area of effect from the studies, and extrapolate to the area for which we are predicting. This is in order to be consistent with the assumptions of the studies whose results we apply. The studies implicitly define a sphere of influence in that the injury risk of road users in one area is influenced only by others within their area and not those in other areas. Therefore, in applying these models, it is consistent to employ a similar or the same area as the sphere of influence, and then scale up to the full area through addition, assuming that risk is constant over the area. Note that to use non-linear coefficients in sizes different from that of the study area defined is a violation of the assumptions inherent in the construction of the regression model.

How does an exponent at a small scale relate to an exponent at a large scale? I.e., what happens to the β coefficients that we learn by fitting a model to small scale if we want to tile small areas, preserving the properties of injury risk encapsulated by α ? To investigate, we define the area of our unit of analysis as A , which corresponds to the typical area size in the study whose results we are applying, and our area of application as nA (e.g. n years). We write the total number of motorists over n years as M and the total number of cyclists C . We approximate the numbers of motorists and cyclists in each unit A as M/n and C/n respectively. We calculate the number of injuries, I_n , as the “sum” over n identical units:

$$I_n = \alpha n \left(\frac{M}{n} \right)^{\beta_1} \left(\frac{C}{n} \right)^{\beta_2}. \quad (9)$$

We can transform the original coefficients, using a function of M (or C) and n , to get the coefficients for the larger scale. Note that M and n describe the relationship between the smaller scale and the larger scale, e.g. when n is close to 1 the coefficients don't change very much.

$$\beta'_1 = \frac{(0.5 - \beta_1) \log(n) + \beta_1 \log(M)}{\log(M)} = \frac{(0.5 - \beta_1) \log(n)}{\log(M)} + \beta_1 \quad (10)$$

We can write this in terms of Equation 2 to show the correspondence, though it obscures the explicit dependence on scale:

$$I_n = \alpha M^{\beta'_1} C^{\beta'_2}. \quad (11)$$

This correction fixes the error of Figure 5, where n is the number of areas (148) and M and C are the total travel in larger geographical regions (rather than local authorities); see Figure 6.

While this allows us to predict injuries for a whole area through consideration of small subunits, it still needs to be extended to model mode shifts and density changes in that larger region. In Section 5 we will join this up with the bigger picture.

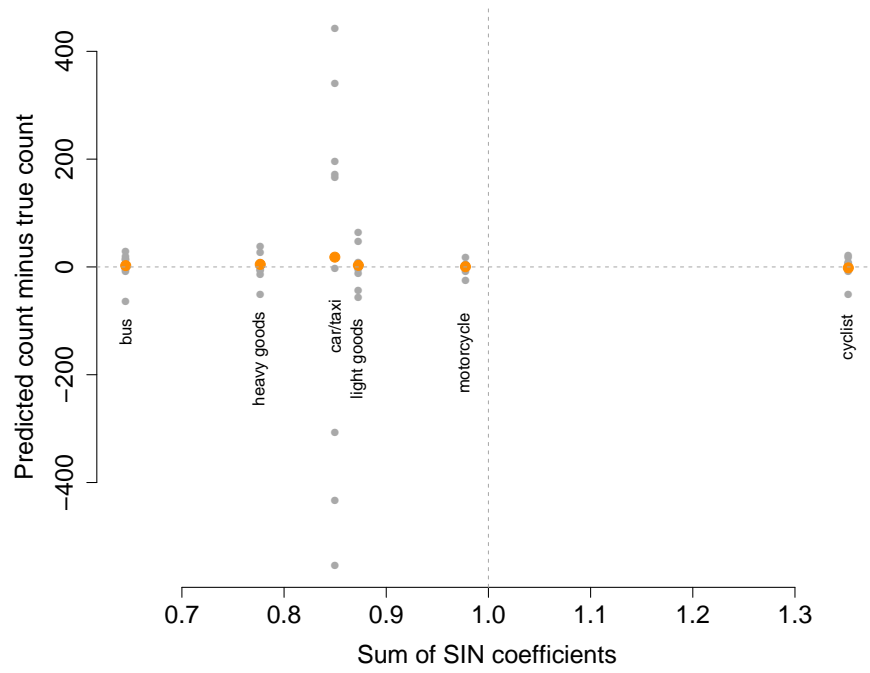


Figure 6: Predictions we make for cyclist KSI counts relative to observed counts, resulting from collisions with other cyclists, motorcyclists, cars, vans, buses and heavy-goods vehicles, as in Figure 5, but with corrected coefficients using Equation 10.

4.2 Mode definitions

Model 2, as is the case for any non-linear model, is not robust to arbitrary redefinition of modes. We ought to be aware of these dynamics when using two-variable models to make predictions in multi-modal settings. Until we know the mechanism for non-linear safety in travel, we won't know how to correctly set up such models.

To illustrate the nature and consequence of the problem, we show how we fit different models if we consider that both “red buses” and “blue buses” are the same “mode”, vs. considering that they are different modes. Let W be the number of red buses, Y the number of blue buses, and $Z = W + Y$ the total number of buses. For collisions with cycles, using equations of the form of Equation 2, we have

$$\lambda_W = \alpha_W W^{\beta_W} C^{\beta_2}, \quad (12)$$

$$\lambda_Y = \alpha_Y Y^{\beta_Y} C^{\beta_2}, \quad (13)$$

$$\lambda_Z = \alpha_Z Z^{\beta_Z} C^{\beta_2}. \quad (14)$$

Let's assume there are three times the number of red buses as blue buses, assume they have the same rate of causing injury to cyclists, and use some simple numbers to illustrate four areas that vary only in size, and not in road-user density or individual-level risks, as given in Table 6.

Table 6: Idealised bus and cyclist injury data.

Observation	C	Buses (Z)	I_Z	Red (W)	I_W	Blue (Y)	I_Y
1	10	80	4	60	3	20	1
2	20	160	8	120	6	40	2
3	30	240	12	180	9	60	3
4	40	320	16	240	12	80	4

As with the illustration in Section 3.1.1, we can fit these models simply, finding in each case that $\beta_i = 0.5 \forall i$. The intercepts, i.e. the base rates, are:

$$\alpha_W = 0.12, \quad (15)$$

$$\alpha_Y = 0.07, \quad (16)$$

$$\alpha_Z = 0.14. \quad (17)$$

It appears that both bus categories are “safer” when considered alone, with blue buses being safer than red buses, despite all bus categorisations posing the same actual danger. This is due only to the non-linear model and the arbitrary definition of one of the covariates.

We arrive at the same conclusion by defining $\exists \alpha : \alpha = \alpha_W = \alpha_Y = \alpha_Z$, reflecting our definition that all bus mode categorisations pose the same risk, and testing whether $\lambda_Z = \lambda_W + \lambda_Y$:

$$\lambda_W + \lambda_Y = \alpha_W W^{\beta_W} C^{\beta_2} + \alpha_Y Y^{\beta_Y} C^{\beta_2}, \quad (18)$$

$$= (\alpha_W W^{\beta_W} + \alpha_Y Y^{\beta_Y}) C^{\beta_2}, \quad (19)$$

$$= (\alpha_W W^{0.5} + \alpha_Y Y^{0.5}) C^{0.5}, \quad (20)$$

$$\stackrel{\exists \alpha}{=} \alpha (W^{0.5} + Y^{0.5}) C^{0.5}, \quad (21)$$

$$\neq \alpha (W + Y)^{0.5} C^{0.5} = \lambda_Z. \quad (22)$$

In order to use non-linear equations for prediction in health-impact modelling in multi-modal settings, it is necessary to provide a solution to this curious problem, or at least a rationale for our choice of approximate solution. The real analogue of this example is the combination and disaggregation of modes in injury modelling. Buses are variously combined with trucks or minibuses; cars, vans and taxis can be one category, two, or three; scooters, motorbikes and three-wheelers can also be considered together or independently.

4.3 Number of modes

Our predictions are highly sensitive to the number of modes we include in the model. Just as we can arbitrarily redefine buses as “red buses” and “blue buses”, if we use Equation 2, we can arbitrarily assume that an injury rate depends on one mode or two. This has profound consequences for our prediction methodology.

To simplify exposition, we assume that injuries modelled as a function of one mode are linear in that mode, and injuries modelled as a function f of two modes are linear in the product of those two modes. More explicitly,

$$\lambda = f(C) = \alpha_c C, \tag{23}$$

$$\lambda = f(M, C) = \alpha_{mc} M^{0.5} C^{0.5}, \tag{24}$$

where we use subscripts c and mc to distinguish the two models. These models are both consistent with the toy example in Table 3. (Note that they are “consistent” with the data in that they describe it. We will go on to look at the problem of interpreting something that is descriptively consistent as something that is a plausible data-generating mechanism.)

The problem arises in our predictive model when we predict injuries based on one mode alone, and keep another constant, i.e. if we want to estimate the causal effect of changing the use of one mode. Suppose we predict the number of injuries to cyclists due to collisions with trucks. We do not change truck distance in our scenario, so a model of the type of Equation 23 seems appropriate, as we are predicting the number of cyclist injuries as a function of cyclists, not of cyclists and trucks. However, they produce very different predictions for the case of doubling cyclists: 2λ injuries predicted from Equation 23, vs. $\sqrt{2}\lambda$ injuries predicted from Equation 24.

What are the consequences of this insight for predictive models where one mode changes not at all, or only a little? Does this bring us back to the problem discussed in Section 3.3? If yes, how ought we rewrite e.g. Equation 10 to account for the varying dependence of the function on the mode arguments? We explore some of these questions in the next section (Section 5). One solution would be to account for scale, which would allow us to include both modes in a consistent manner, and which we return to in Section 6.

5 Simulation and motivation for an alternative model

Here we look explicitly at city-level models, that is, studies where the unit of observation is the city (rather than small localities such as junctions) in order to understand how safety scales as a function of road use. As at other scales, coefficients $\beta_1, \beta_2 < 1$ have conventionally been understood to correspond to safety in numbers and, again, coefficients learnt in such studies have been suggested as parameters to be used in predictions of injury numbers in mode-shift scenarios i.e. causal effects of mode usage.

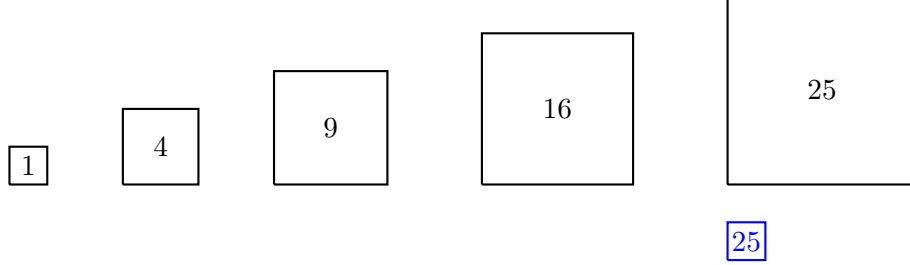


Figure 7: Schematic representation of an inter-city study, demonstrating scaling across cities. If we fit the model in Equation 2 using these cities, and then make a prediction for “City 1” with 25 times the travel, then the best estimate we make will correspond to the observation for “City 25” (black). However, when modelling mode-shift scenarios, what we actually want to predict is represented in blue: the total travel of “City 25”, taking place in a city of the size of “City 1”.

Figure 7 demonstrates the fundamental problem in using cities that vary across scales to inform models intended for use in predicting for a city that will change in number but not scale, i.e. it will change in density. In this section we develop the ideas summarised in Figure 7 through a simulation study based on a simple individual-level physical mechanism for generating collisions between different parties, and injuries.

5.1 Simulation model

We create a simulator with three variables: the number of cyclists, the number of motorists, and the dimension of the space, which we equate to the size of the city. Each cyclist and each motorist is a five-pixel by five-pixel square. At each time step each body moves one body-width. We simulate 500 time steps. Each body undertakes a biased random walk: that is, with probability 5/6 it continues in the same direction, and with probability 1/6 it chooses a direction randomly from among the directions available to it (including the same direction). Moves that would take a body out of the frame wrap around to the other side of the frame, i.e. emulating balanced migration into and out of the area. The frame size is some multiple of the step size that we vary as an input parameter.

We count the number of collisions between cyclists and motorists, which is defined as any overlap in pixels between a cyclist and a motorist. In the style of PacMan, upon collision the cyclist disappears and reappears randomly in the next time step. The cyclist disappears immediately, so a collision involving two motorists and one cyclist will be counted as one event. There is no “safety in numbers” in this model: a cyclist’s risk of collision is independent of the number of other cyclists in the simulation.

5.2 Simulation study

In order to emulate inter-city regression studies, we simulate 50 times frames of different sizes and constant density. We define our vector of sizes to be the integers $x = \{5, \dots, 14\}$. The numbers of cyclists and motorists are Poisson-distributed random variables with mean x^2 , and the dimensions of the frames are $20 \times 5x$ pixels by $20 \times 5x$ pixels. There are 500 simulations in total: ten sizes with fifty repetitions of each. The results of the simulations are summarised in Figure 8. We fit a Poisson regression model to the resulting number of collisions to learn the parameters β_1 and β_2 , finding

$\beta_1 \approx \beta_2 \approx 0.5$ as expected, since the density remains constant, but the area sizes vary. We then use (a) the model to make new predictions and (b) the simulator to test them.

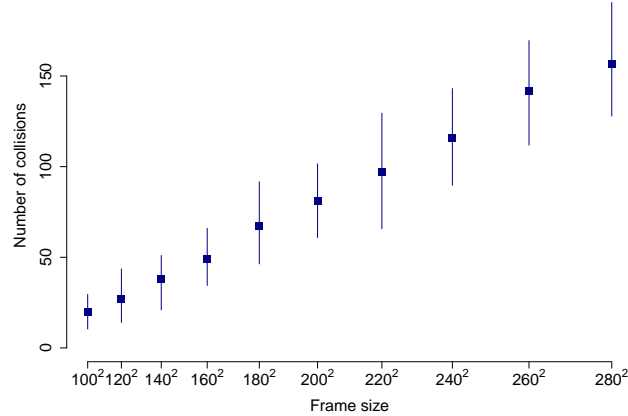


Figure 8: Simulation study. The number of collisions is a function of frame size, number of cyclists and number of motorists. Each bar shows 90% of the range of 50 simulations. In every simulation the density of cyclists and the density of motorists is the same.

5.3 Comparison of model predictions to simulations

First, we consider every frame size with 100 motorists and 100 cyclists. This was one of our simulation inputs for generating the simulated data: frame size 200². Because our model is a function of road-user numbers alone, we make the same prediction for every frame size: 80 collisions. In Figure 9 we plot this against the simulated number over 50 repetitions for each frame size, demonstrating the extent of the bias resulting from having left frame size out of the regression model. Note that only for the frame size corresponding to mode numbers of 100 does the 90% confidence range overlap the predicted value.

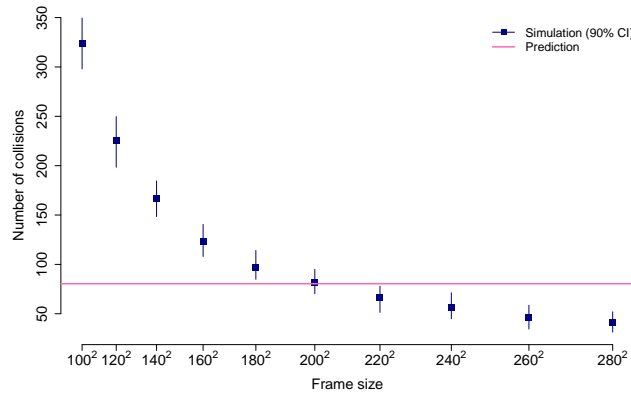


Figure 9: Simulated vs. predicted results for 100 motorists and 100 cyclists with varying frame sizes. The predicted value is 80 for all frame sizes.

Second, using the same set up, we consider that the number of motorists is constant at 100 in the corresponding frame size of 200², and we vary the number of cyclists. Again, we predict the number of collisions using our model, which this time varies with cyclist number. Note in Figure 10 that the prediction aligns with the simulation when the cyclist value corresponds to the frame size and density of the predictive model. When this number is exceeded, the estimate is too low, and at a lower density, the prediction is too high. We should keep this picture in mind when we

consider e.g. cyclist injuries when cyclist numbers are changing but the other mode, such as truck, is not.

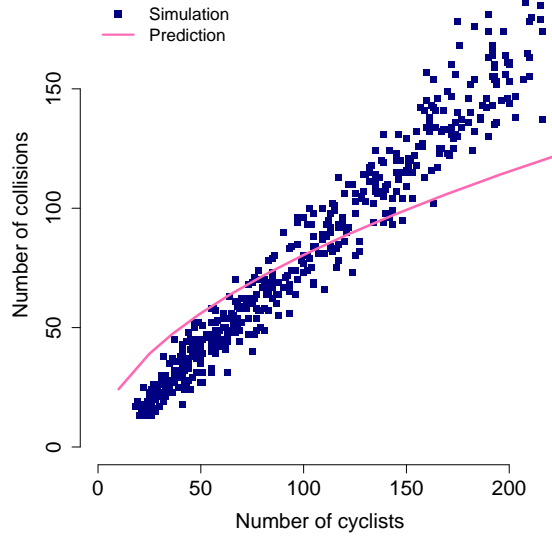


Figure 10: We fix motorists at 100 and frame size at 200^2 , and vary the number of cyclists as Poisson random variables with means from 25 to 196. We simulate and predict using the model the number of collisions. This represents the type of error we might make if we predict for one mode that varies when another stays constant.

Third, we consider that the number of cyclists is constant at 150 in a frame size of 200^2 , and we vary the number of motorists. Again, we predict the number of collisions using our model, which this time varies with motorist number. Note in Figure 11 that the prediction aligns with the simulation when the motorist value corresponds to the difference between the motorist value corresponding to the frame size and density of the predictive model, and the change in cyclist number from its corresponding value. (Precisely: the errors cancel out when the number of new motorists = the number of old motorists \times the number of old cyclists / the number of new cyclists = $100 \times 100 / 150 = 66.7$.) When this number is exceeded, the estimate is too low, and at a lower density, the prediction is too high. We should keep this picture in mind when we consider e.g. that we have six modes each with 100 road users, and we predict for a scenario in which cyclists increase to 150 and all other modes decrease to 90.

5.4 Simulated safety in numbers

We repeat the simulation study, this time imposing a safety-in-numbers effect. We specify that the probability of collision is represented by some number $p = p(C)$, which is a function of the number of cyclists C . We recreate Figure 8 with $p = C^{-0.25}$. In terms of “safety in numbers”, this corresponds to a raw exponent of 0.75. The probability p doesn’t depend on M , so its exponent is 1.

We see the effect of cyclist number on collision number in Figure 12. Note the non-linear gradient, compared to Figure 8. For these data, we infer for a model of the form of Equation 2 exponents $\beta_1 + \beta_2 \approx 0.75$, which is equal to the sum of the raw exponents input into the simulation (defined above) minus 1: $\beta_1 + \beta_2 \approx 0.75 = 1 + 0.75 - 1$.

Next, we simulate the number of collisions for this model with doubled numbers of cyclists and motorists and varying frame sizes. The results are shown in Figure 13 along with the prediction from the exponents inferred before ($\beta_1 + \beta_2 \approx 0.75$) from Figure 12.

Concatenating simulated data in Figures 12 and 13, we fit a model with exponents $\beta_1 + \beta_2 \approx 1$, recovering the tiling parameters discussed previously, showing that the model is fitting across scales – hence we refer to them as scaling exponents. When we control for frame size (i.e. we specify a

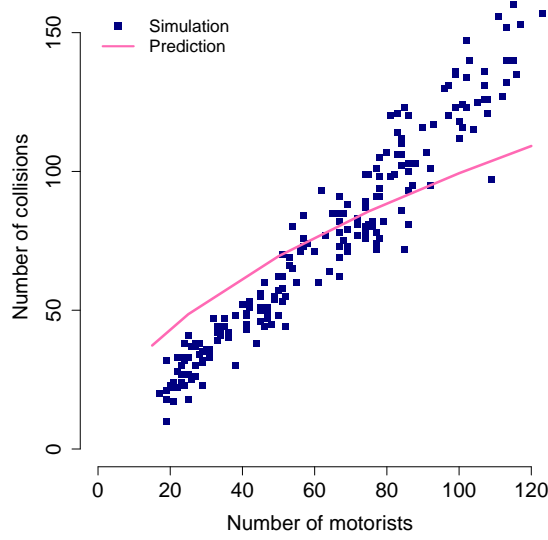


Figure 11: We fix cyclists at 150 and frame size at 200^2 , and vary the number of motorists as Poisson random variables with means from 25 to 100. We simulate and predict using the model the number of collisions. This represents the type of error we might make if we predict for one mode that varies as a result of reallocation of other mode types. Note the deviation between model and prediction, which is minimised close to the point of constant density, i.e. that the increase in cyclists is close to the decrease in motorists.

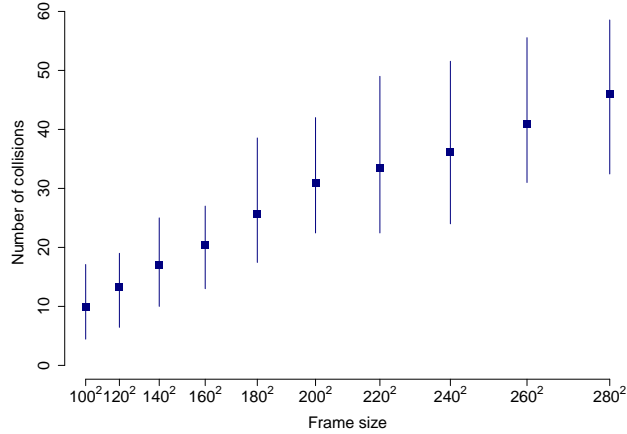


Figure 12: Simulation study. The number of collisions is a function of frame size, number of cyclists and number of motorists. Each bar shows 90% of the range of 50 simulations. In every simulation the density of cyclists and the density of motorists is the same. Collisions are a function of the number of cyclists and occur with probability $p = C^{-0.25}$.

different model (Equation 25), which eliminates confounding), we isolate the density effect: we will refer to these exponents (those belonging to the density model, which we formally define later in Section 6) as δ_1 and δ_2 , analogous to β_1 and β_2 . We find that $\delta_1 + \delta_2 \approx 1.78$, close to the original exponents.

We repeat the whole process for $p = C^{-0.5}$, i.e. exponents 0.5 and 1. Analogously for Figure 12 we find $\beta_1 + \beta_2 \approx 0.5$, the sum of the original exponents minus one. Again we simulate data for doubled cyclists and motorists which, concatenated to the original data, yield scaling exponents $\beta_1 + \beta_2 \approx 1$ and density exponents $\delta_1 + \delta_2 \approx 1.5$ when size is accounted for.

We are trying to build up a picture of how Equation 2 behaves across scales and densities. It

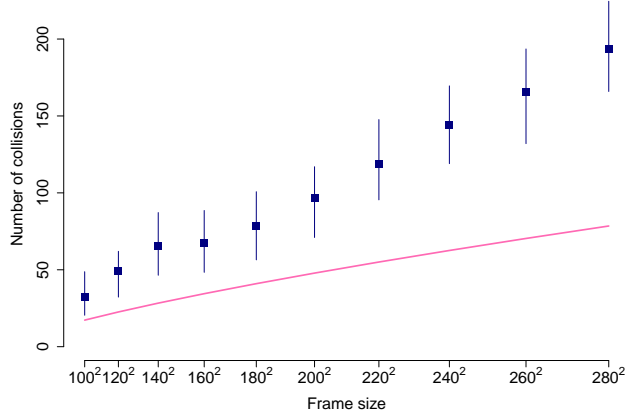


Figure 13: As in Figure 12, with double the number of motorists and cyclists. The prediction generated from data in Figure 12 is shown in pink.

seems from Figures 9, 10, 11 and 13 that the predictive equation we are aiming for requires density exponents; that studies across scales return scaling exponents; and that, for two modes, the sum of the density exponents is one more than the sum of the scaling exponents: $\delta_1 + \delta_2 \approx \beta_1 + \beta_2 + 1$. We will explore in the following Section if this observation generalises, and what its mechanistic basis might be.

5.5 Discussion

We do not simulate in order to suggest that we have understood and can recapitulate the mechanism. Rather, we use a simple simulation model in order to assess the effect of model (mis)specification on model predictions and infer simple principles, for example the relationship between size and density. If we can make a general statement about what happens in this simple simulation, in which we have full control of the workings, we have the opportunity to illuminate possible factors influencing observational inferences such as those in safety-in-numbers studies.

There are some key differences between our simulation study and real data: we have no spatial or temporal factors governing the rate of collision, and no relationship between e.g. density and speed. One respect in which our simulation differs from the England study of Section 3.3.2 is that we consider here constant density of road users across scales, whereas, on average, the areas in England decline in density as size increases. In contrast, in general, it is posited that, globally, as city size increases, so does density (see schematic of the model space, Figure 14). However, we could tailor our simulation to match in some way a phenomenon we are interested to capture.³

The point here is not how much our simulation set up resembles what one imagines happens in cities. The point is that we can contrive data using a transparent physical mechanism from which we can learn coefficients that fit the general trend $\beta_1 \approx \beta_2 \approx 0.5$ using Equation 2, and we can test this model against data simulated from the same original source. This means we can test the range of applicability of the model. We simulated data that exhibits no safety in numbers, and we inferred a safety-in-numbers effect using Equation 2. To correct for this bias, we used a size-adjusted model that accounts explicitly for scale, and which we will expand upon in the next Section.

5.6 Conclusion

The conclusion of this test is that a predictive model that does not take into account factors of scale fails to predict outside its training space. We can identify the direction of the bias: uncaptured

³For example, we could introduce density-dependent speeds. Concretely, in a fixed space, doubling cyclists and doubling motorists increases density by a factor of four. If we then divide every mode speed by that value - four - we undo exactly the effect of the increased density.

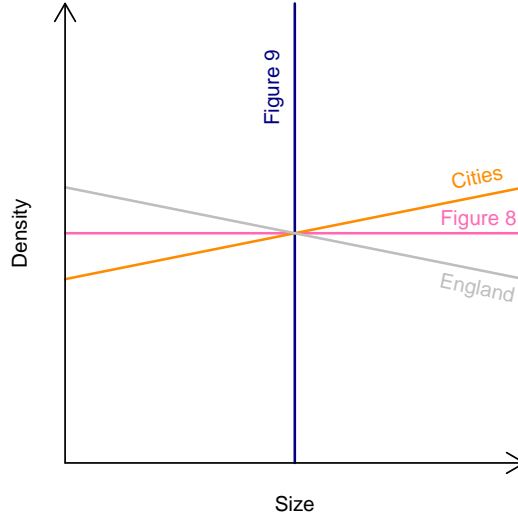


Figure 14: A depiction of the model space covered by the simulation model. We mark out in pink and navy blue a line of constant density and a line of constant size through the model space, depicted in Figures 8 and 9, respectively. The study of areas in England might be represented in this model as occupying the space of increasing size and decreasing density (grey), and expect that a study of many cities will occupy the space of increasing size and increasing density (orange). (NB: these are schematic.)

increases in density lead to underprediction of collisions, and uncaptured decreases in density lead to overprediction. This calls into question the assumption that density-independent regression models can be used to predict the number of collisions, or injuries, that will occur in mode-shift scenarios.

A heuristic for correction is offered (formalised in Section 6): for our simulation, scaling exponents of 0.5 corresponded to density exponents of 1, with two modes simulated. A starting point for application of inter-city studies is then the translation of a pair of scaling exponents β_1, β_2 to a pair of density exponents $\beta'_1 + \beta'_2 = \beta_1 + \beta_2 + 1$. This would correct the errors seen in Figures 9, 10 and 11. Then our null hypotheses, assuming road-user risk per individual road user is independent of all other road users of that type, are, for models of the form of Equation 2, that $\beta_1 + \beta_2 = 1$ for constant density and varying scale, and $\beta_1 + \beta_2 = 2$ for constant scale and varying density.

6 An alternative model

Recall that our null hypotheses are, for models of the form of Equation 2, that

- (1) $\beta_1 + \beta_2 = 1$ for models fitted to datasets in which there is constant density and varying scale, and
- (2) $\beta_1 + \beta_2 = 2$ for models fitted to datasets in which there is constant scale and varying density.

Testing these hypotheses would allow us to learn about the effect of road-user number on injury risks. However, as it's unlikely that we will find datasets that have either constant size or constant density, we create a model that brings both together, so that scale/density is accounted for in the model. The model we write to capture both those special cases is parametrised by exponents δ_1 and δ_2 and also by size; i.e. (1) and (2) are “special cases” and complementary presentations of the overarching model. Here, we present the model used in place of Equation 2 to correctly capture the dynamics in the simulations for general use.

6.1 The model

We assume that there is a constant number of road per unit in order to simplify the model and its derivation, i.e., we consider that there are n units, with $c = C/n$ cyclists per unit and $m = M/n$ motorists per unit. A more precise model would accommodate heterogeneity in road-use density. Then our model is written:

$$I_n = \frac{\alpha}{n} (nm)^{\delta_1} (nc)^{\delta_2} = \frac{\alpha}{n} M^{\delta_1} C^{\delta_2} \quad (25)$$

where n is the number of spatial units.⁴ Then, in a scenario, we can distinguish between a doubling of density,

$$I_n^{(\text{double density})} = \frac{\alpha}{n} (2M)^{\delta_1} (2C)^{\delta_2} \quad (26)$$

$$= 2^{\delta_1 + \delta_2} \frac{\alpha}{n} M^{\delta_1} C^{\delta_2} \quad (27)$$

$$= 2^{\delta_1 + \delta_2} I_n, \quad (28)$$

and a doubling of size, i.e. we double the number of areas n , which has the consequence of doubling M and C :

$$I_n^{(\text{double size})} = \frac{\alpha}{2n} (2M)^{\delta_1} (2C)^{\delta_2} \quad (29)$$

$$= 2^{\delta_1 + \delta_2 - 1} \frac{\alpha}{n} M^{\delta_1} C^{\delta_2} \quad (30)$$

$$= 2^{\delta_1 + \delta_2 - 1} I_n. \quad (31)$$

When density doubles, because M and C double and the size remains the same, we expect injury number to multiply by $2^{\delta_1 + \delta_2}$. When density stays the same, so that when M and C double, so too does the size, we expect the injury number to multiply by $2^{\delta_1 + \delta_2 - 1}$.

With this formulation we can also derive the relationship with β_1 and β_2 , assuming observations for a size of 1 and a size of n . From Equation 25 we have

$$\frac{I_n}{I_1} = \frac{\frac{\alpha}{n} (nm)^{\delta_1} (nc)^{\delta_2}}{\alpha m^{\delta_1} c^{\delta_2}} \quad (32)$$

$$= n^{\delta_1 + \delta_2 - 1}. \quad (33)$$

⁴Note that this formulation is equivalent to Equation 9 for $\beta_1 = \delta_1$ and $\beta_2 = \delta_2$ and $\delta_1 + \delta_2 = 2$. Equation 9 computes the number of injuries in a whole space as the sum of the number of injuries in its subspaces, and has the special property that for $\beta_1 + \beta_2 = 1$ the number of injuries is independent of the size of the space given M and C . The formulation of Equation 25 does not admit this independence.

We can derive the equivalent relation using Equation 2:

$$\frac{I_n}{I_1} = \frac{\alpha(nm)^{\beta_1}(nc)^{\beta_2}}{\alpha m^{\beta_1} c^{\beta_2}} \quad (34)$$

$$= n^{\beta_1 + \beta_2}. \quad (35)$$

Thus we see that for these two models to be consistent, we require $\beta_1 + \beta_2 + 1 = \delta_1 + \delta_2$.

6.2 Discussion

We propose the model of Equation 25 as a model that is consistent with our observations and our insights. There will be other models that also fulfill those criteria. There will be other models that contain ours within them as a subset or a special case. An example would be an extension that takes account also of speed. Whatever model we use, however, ought to be consistent with the observations made so far, and this is one such model.

In terms of application, we know how to transform $\beta_1 + \beta_2$, but we don't know how to transform β_1 and β_2 separately. We could define $\delta_i = \beta_i + 0.5$, adding an equal amount to each variable. This assumes that a cross-scale application of the model in Equation 2 is accurately capturing the contributions of the individual modes to risk. However, inference of β_1 and β_2 is likely confounded by space and/or time, so we might instead choose $\delta_1 = \delta_2 = (\beta_1 + \beta_2 + 1)/2$, an equitable solution, which assumes that inference using Equation 2 has identified the correct sum $(\beta_1 + \beta_2)$ but not their individual contributions. These options should be tested, particularly for cases where one mode vastly outnumbers the other. For models expressing uncertainty, distributions can be assigned to these parameters and a sensitivity analysis conducted to ascertain the impact of the values on the outcome.

Finally, the problem identified in Section 4.3 remains for $\delta_1 \neq 1$ where δ_1 is applied to a mode that is arbitrarily redefined. A quick fix would be to define $\delta_1 = 1$ for this mode, or group of modes, and $\delta_2 = \beta_1 + \beta_2$. Of course this won't help when both modes are arbitrarily defined. In this case, one can re-parametrise the model, or make multiple assessments (one per mode definition) as a sensitivity analysis.

7 A size-adjusted study of the England data

We apply the previously presented model (Equation 25) to the England data, using the distance travelled as covariates and the total road length for each area as offsets, assuming a Poisson distribution for the counts. We use the 148 areas of England and consider only A, B and minor roads. We consider urban and rural areas both separately and together, defining urban areas as those with at least 98% of their population registered as living in a city, town or minor conurbation in the 2011 census (2079 data points; 2889 for rural). We use the software Stan in R with default (uniform) priors to test the hypothesis $\delta_1 + \delta_2 = 2$ by evaluating the posterior probability that $\delta_1 + \delta_2 < 2$ and inspecting its distribution, e.g. whether it is concentrated around a value that is practically different from 2.

7.1 Results

In Table 7 we present the 95% credible intervals for the sum $\delta_1 + \delta_2 = 2$. Our data and model give good evidence against the hypothesis of linearity ($\delta_1 + \delta_2 = 2$) for all injuries and KSI in urban areas, but not in rural areas. The evidence against this hypothesis for fatalities in all areas is weak. The ranges are much larger for fatalities, presumably because the signal is weaker as there are fewer events and more zeros. It might be that with more data we would be able to make stronger statements with more confidence about the data in relation to the null hypothesis.

We note that the sums are about 1 more than the sums $\beta_1 + \beta_2$ that we learn with the same data and a simpler model that considers numbers alone and not road length (Table 2 for “All areas”: 0.67 for all injuries, 0.79 for KSI, and 0.97 for fatalities). (In fact, the difference is slightly greater than 1, as the areas in England have a slightly negative correlation between size and cumulative cycling, and hence greater adjustment is required, depicted as the gap between the grey and navy lines in the schematic in Figure 14, which exceeds 90° .)

Table 7: 95% credible intervals for $\delta_1 + \delta_2$ and probability < 2 .

	All injuries		KSI		Fatalities	
	95% CI	Probability	95% CI	Probability	95% CI	Probability
All areas	1.79–1.80	1.00	1.90–1.93	1.00	2.00–2.18	0.02
Urban areas	1.62–1.65	1.00	1.74–1.83	1.00	1.42–2.08	0.93
Rural areas	1.92–1.93	1.00	1.96–2.01	0.92	1.90–2.14	0.37

In Tables 8 and 9 we present the 95% credible intervals for the coefficients for cycle and car travel, and in Table 10 the intervals for the intercepts. Note the negative correlations between coefficients: for “All areas” in Tables 8 and 9, as casualty severity increases, more of the coefficients’ sum is attributed to car at the expense of cycle. Similarly, across all urbanicity levels, there is a negative correlation between the sums in Table 7 and Table 10, suggesting a tradeoff between base rate and safety effect.

Table 8: 95% credible intervals for δ_2 (cycle travel) and probability < 1 .

	All injuries		KSI		Fatalities	
	95% CI	Probability	95% CI	Probability	95% CI	Probability
All areas	0.59–0.61	1.00	0.52–0.57	1.00	-0.01–0.30	1.00
Urban areas	0.68–0.71	1.00	0.64–0.72	1.00	-0.00–0.49	1.00
Rural areas	0.19–0.22	1.00	0.18–0.26	1.00	-0.16–0.28	1.00

Finally, we present the coefficients we estimate when we combine all datasets, including the new levels as factor predictors in a multivariable regression model (Table 11). Using factors for casualty severity and urbanisation, we end up with coefficients that resemble those fit to the dataset with all injuries and pooled urbanisation.

Table 9: 95% credible intervals for δ_1 (car travel) and probability < 1 .

	All injuries		KSI		Fatalities	
	95% CI	Probability	95% CI	Probability	95% CI	Probability
All areas	1.18–1.20	0.00	1.35–1.39	0.00	1.79–2.10	0.00
Urban areas	0.93–0.95	1.00	1.06–1.14	0.00	1.20–1.80	0.00
Rural areas	1.70–1.73	0.00	1.72–1.81	0.00	1.72–2.23	0.00

Table 10: 95% credible intervals for α (the intercept).

	All injuries	KSI	Fatalities
All areas	-22.97–22.76	-27.87–27.27	-38.24–34.59
Urban areas	-19.65–19.06	-25.29–23.55	-35.57–23.69
Rural areas	-27.52–27.16	-30.91–29.94	-37.93–32.52

7.2 Discussion

We present the test of $\delta_1 + \delta_2 = 2$ as the test for “size-adjusted safety in numbers”, which can be done in the Bayesian framework with Stan. Using glm in R, we can test $\delta_1 = 1$ and $\delta_2 = 1$ (i.e., test per capita as in Shalizi (2011)) using the offset function, but we cannot test their sum.

Both in the glm framework and with Stan, it is unclear to what extent we can interpret the values for δ_1 and δ_2 . Were we to use data in which there is no correlation at all between cycling volume and car volume, they will be uniquely identifiable parameters. On the other hand, if in our data cycling and car volumes were perfectly correlated, we could not separate out separate values at all. In reality, cycling and car volume are somewhat correlated, and therefore the values for δ_1 and δ_2 are somewhat entangled. Therefore, the values presented in the tables are unlikely to be representative of causal effects. It is possible that there is more safety in numbers for cyclists in rural areas, and a corresponding danger in numbers from cars in rural areas. It is also possible, given our model and data, that there is sharing of the coefficients when the model is fit, so that both values should be closer to 0.95.

There appear to be correlations between coefficients across the models. What does it mean that, for “All areas” in Tables 8 and 9, as casualty severity increases, more of the coefficients’ sum is attributed to car at the expense of cycle? It seems likely that there is a pattern underlying the data, which is consistent across severity levels, and that the model is unable to describe. It seems less likely that casualty severity impacts on the dynamics of the non-linear relationship between distance travelled and road-traffic collision rates. There is also a negative correlation between base rate and safety effect in Tables 10 and 7: do places that are more safe have less of a safety effect? Or is there some transference between the parameters?

This analysis is comparable to the population-adjusted model of Aldred et al. (2017). In that study, there are three covariates: cycle commuters, motor vehicle volume, and population. Here, we have two covariates, cycle distance and motor distance, and one offset, the total length of A, B and minor roads.

The number of cycling commuters is similar to estimated cycle distance; motor vehicle volume and motor distance are measuring the same thing; and population and road length are correlated, and are both a proxy for an area’s size. The key difference between the two models is the treatment of the city proxy as a covariate (population) and as an offset (road length), which is equivalent to a covariate with a fixed coefficient, which we set to -1.

The sums of the means for the covariate coefficients of the population-adjusted model of Aldred et al. (2017) are 0.98, 1.06, 1.02, and 0.99. We cannot test a hypothesis in this regression framework but it’s worth noting the proximity to our corresponding estimates, $\delta_1 + \delta_2 - 1$, which are consistently just less than 1.

That there is a systematic difference between “urban” and “rural” areas in terms of the coef-

Table 11: 95% credible intervals the δ coefficients using a factor covariate for severity and probabilities of the values being less than 1 and their sum being less than 2.

	δ_2 (cycle travel)		δ_1 (car travel)		$\delta_1 + \delta_2$	
	95% CI	$P(\delta_2 < 1)$	95% CI	$P(\delta_1 < 1)$	95% CI	$P(\delta_1 + \delta_2 < 2)$
All areas	0.59–0.61	1.00	1.18–1.20	0.00	1.79–1.80	1.00
Urban areas	0.68–0.71	1.00	0.93–0.95	1.00	1.62–1.65	1.00
Rural areas	0.19–0.23	1.00	1.70–1.73	0.00	1.92–1.93	1.00
With factor	0.58–0.60	1.00	1.23–1.25	0.00	1.83–1.84	1.00

ficients fit by our model is indicative that the model does not capture the whole effect of travel density on injury rates. This, and the spurious correlations in coefficients, show there are features of the data not explained by the model of Equation 25.⁵

Mis-specification of any of (a) the relationship in the model between size and rate, (b) the probabilistic description of the error term, (c) the component contributions of the two mode distances, and (d) quantification of size through road length might contribute to the failure to explain the data, and all would benefit from further consideration and testing. All of these things can be improved upon, incrementally, as we have taken an incremental step from Equation 2 to a size-based model in Equation 25. However, taken together, they highlight that it has not been shown that a model with the fundamental form of Equation 2 might be capable of answering the question of safety linearity.

⁵The difference does, however, open the possibility that there is “nonlinearity in the nonlinearity” – i.e., that the effect is greater (and the exponents smaller) at higher densities. This is an avenue that could be explored as a link between city-scale and small-scale studies, which typically consider only high-density areas of a city.

8 Conclusions

For studies that scale across sizes, $\beta_1 + \beta_2 = 1$ for Equation 2 represents linearity in numbers. That the same (or similar) coefficients are observed across scales is perhaps the biggest indication that whatever is being captured is not an effect of cyclists conferring protection to other nearby cyclists.

The missing component in Equation 2 is size. Omission of size results in misleading interpretations and predictive models unsuitable for fixed settings. Therefore, we recommend departing from this model and developing new expressions, which include size explicitly. We propose, in the first instance, a very simple adjustment to Equation 2 in Equation 25, in which we include size as defined by total road length.

We recognise the challenge of specifying testable hypotheses in this setting. We are trying to formulate the hypothesis that the risk to an individual road user colliding with another road user is independent of the number of other road users of their type, and linear in the number of road users of the other type, within a particular space. Then, for each cyclist, the risk would be αm . For c cyclists, the expected number of injuries would be αmc . For n spatial units, the expected number of injuries would be $n\alpha mc$.

Ideally we would test per capita rates for δ_1 and δ_2 in Equation 25, but, as the model stands, they cannot be confidently identified, so we test instead the null hypothesis $\delta_1 + \delta_2 = 2$ in Stan. With this model, hypothesis, and the data for England, we found a size-adjusted safety-in-numbers effect for all subsets of the dataset, with the exception of the set of fatalities in rural areas, whose null hypothesis we do not reject given the data we have.

The model we present is best described as “size-adjusted safety in numbers”; we have not developed or tested a model of “safety in density”, which might be interesting and relevant to explore. It would allow for heterogeneity in density over space and time and exploration of their impacts on results. That will be particularly important for small-scale studies, e.g. of junctions, and could be addressed in the first instance through simulations such as those presented in Section 5. In addition to heterogeneity in density over time and space, an improvement to the model in Equation 25 would be inclusion of mode speeds.

In terms of hypothesis testing, we aim to engage a wider audience and share data in order to find alternative ways to understand and formulate the problem. We recognise a correspondence between our model and city-level metrics that were once claimed to exhibit power-law scaling properties (Leitão et al., 2016) (see Appendix C). Input from these statistical modellers might greatly benefit progress in this topic. In addition, there might be parallels with the practice of discretisation of space (and time) in a “contact matrix” to describe interactions across partitions in infectious-disease modelling (Birrell et al., 2011). We have identified some areas for further testing or development of the model described by Equation 25, such as specification of the error term, use of density proper rather than size adjustment, the expression of the mode distances, how to quantify size, and the problem of mode (dis)aggregation. We would welcome development of other models, as well as methods for assessment and hypothesis testing.

Complementing data-driven analyses, simulation models can be used to develop relational models. These can test implications of comprehensive mechanistic models of injuries as a function of space and its occupancy. Simulation and theory provided us with a number of null hypotheses, which (a) give us an objective to test and reject, and (b) provide a justified basis for prediction. In this work, the theory and simulation led to the following hypotheses, which connect Equation 2 to Equation 25:

- (1) We can estimate the city-scaling exponent from small-scale exponents by scaling β coefficients in such a way that we approximate the process of “multiplying then summing” rather than “summing then multiplying” (Section 4.1).
- (2) City-level size-scaling exponents $\beta_1 = \beta_2 = 0.5$ correspond to linearity (Section 3.3).
- (3) City-level size-scaling exponents $\beta_1 = \beta_2 = 0.5$ (Equation 2) correspond to city-level density-scaling exponents $\delta_1 = \delta_2 = 1$ (Equation 25), where mode speeds are assumed independent of density.
- (4) To predict the consequence of a change in density using city-level size-scaling exponent β , we can use density exponents $\delta_1 + \delta_2 = 1 + \beta_1 + \beta_2$, where there are two modes involved (Section

5.4).

These hypotheses were consistent with the results of the England data, in that the coefficients δ_1 and δ_2 fit with the model of Equation 25 had sum approximately one greater than that of β_1 and β_2 fit with the model of Equation 2. We therefore propose this framework as a first amendment for making predictions of road-injury burden in mode-shift scenarios (question 2 in Table 1). With reference to all those initial objectives, we hope to have offered a new perspective, and to have generated more questions and pointed to new lines of inquiry. Future avenues might include:

- (1) How studies into whether a scaling is non-linear can be applied to this setting (e.g. Leitão et al. (2016)).
- (2) How to explain the relationship between small-scale dynamics and city-level dynamics.
- (3) Inclusion of heterogeneity in density.
- (4) How to include or assess temporal dynamics.
- (5) How to include road length and speed as variables.
- (6) Different types of model, such as using mode shares as predictors.

A Supplementary figures

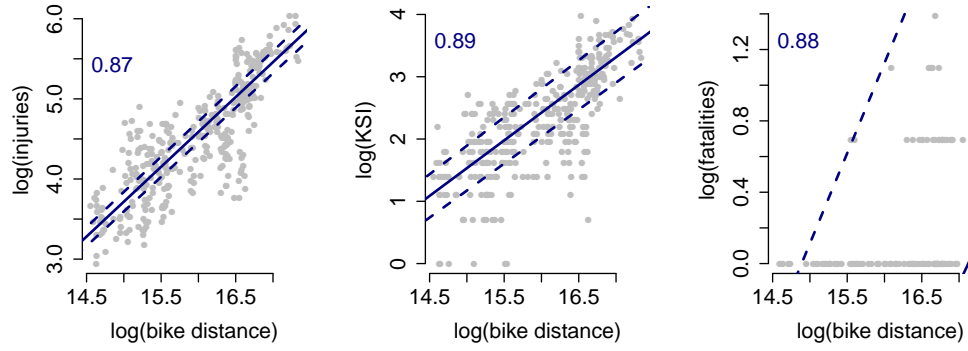


Figure 15: As in Figure 1, with only the London boroughs.

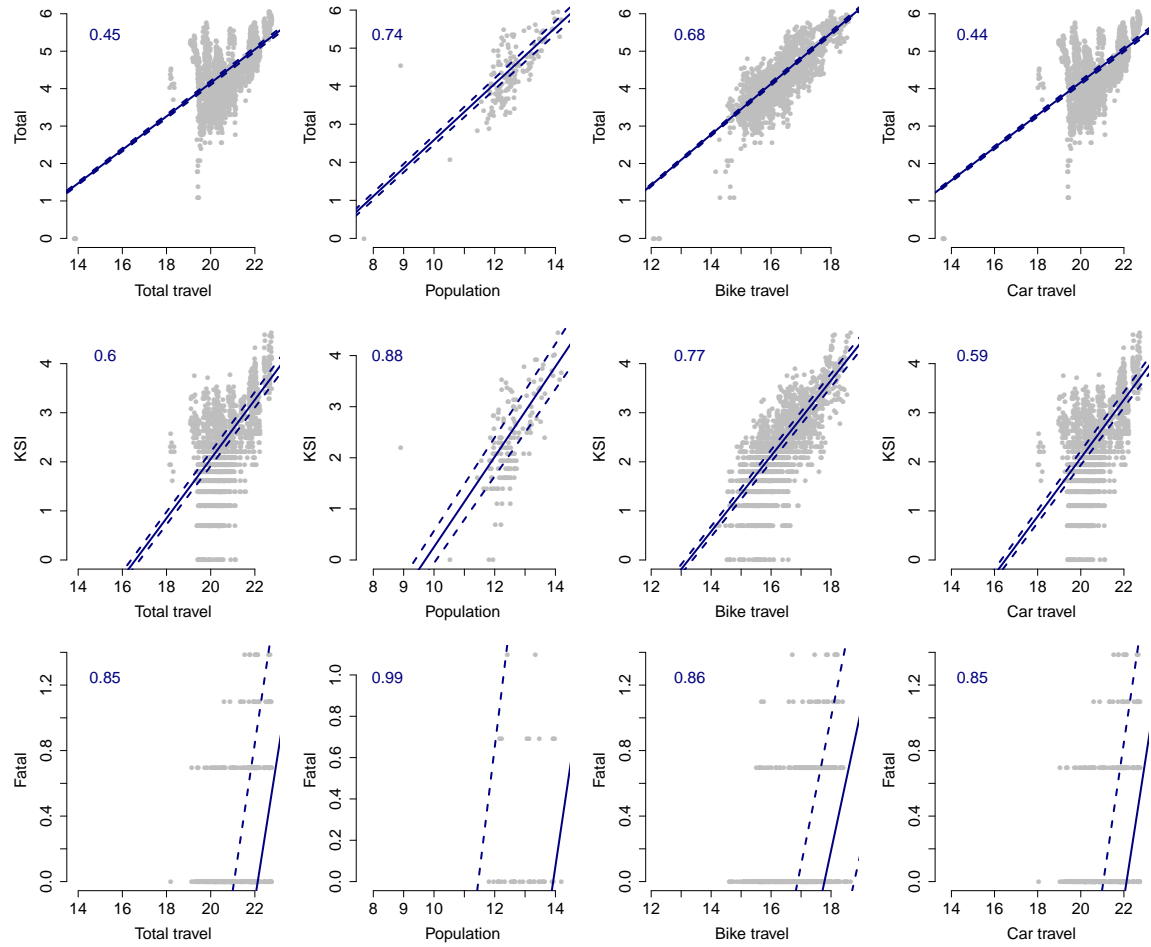


Figure 16: As in Figure 1, showing different relationships and the β coefficients they generate.

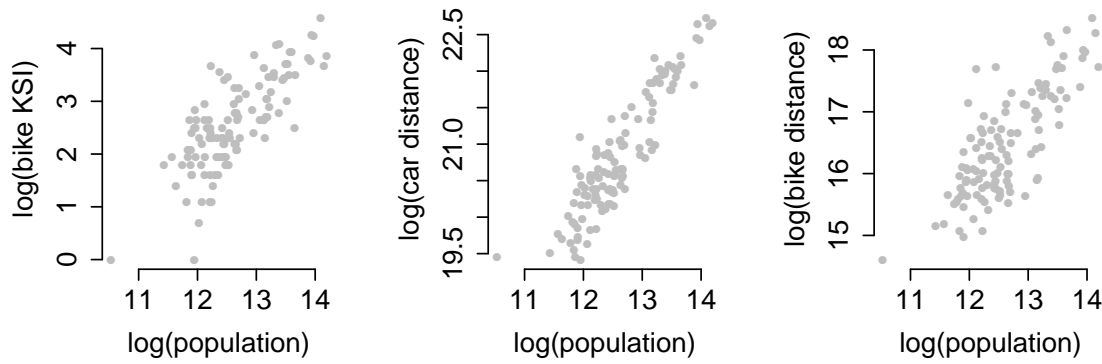


Figure 17: Linear relationships between population (N) and cyclist KSI, between population and car distance, and between population and cyclist distance in English counties.

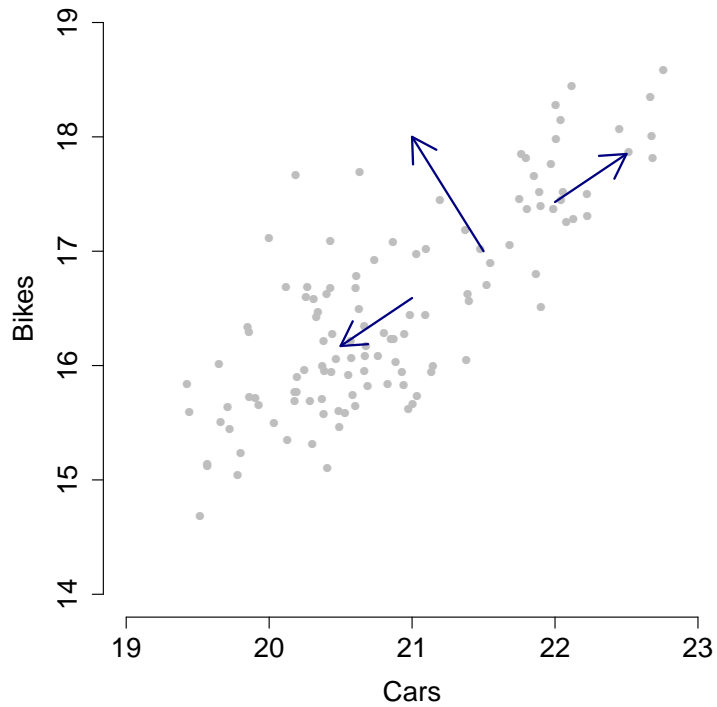


Figure 18: Where we have two colinear predictors, it seems reasonable to make predictions in spaces that align with their colinearity. However, does it make sense to make predictions in spaces orthogonal to the colinearity?

B Worked example: implementation in ITHIM-R

We consider the setting of Accra, for which we have a list of recorded fatalities over multiple years. Each record contains the following information: the year, the mode of the casualty, the mode of the other party, the age of the casualty, and the gender of the casualty. In addition, we have a travel survey, from which we learn total travel by each mode (and by demographic group, which we omit for now, for simplicity).

B.1 Constructing the model

We fit the observed data (the number of injuries, I) to an equation of the form

$$I \sim \text{Poisson}(\lambda), \quad (36)$$

$$\lambda = \alpha M^{\beta_1} C^{\beta_2} \exp \left(\sum_{i=3}^P X_i \beta_i \right) \quad (37)$$

with α a fixed intercept, C and M the distances travelled by cyclists and cars, respectively, based on the travel survey, and X the model matrix built from all the covariates (here, we consider only the two modes; gender and age of the casualty are omitted for simplicity). We do not use the “year” covariate but instead suppose that we have multiple observations for a single “year” (i.e. we reuse the distance data). Finally, the coefficients to fit using `glm` are α and β_i for $i \geq 3$, and we supply β_1 and β_2 as fixed parameters so that $M^{\beta_1} C^{\beta_2}$ is our offset.

Note that there are many combinations of modes, so this model is linked via the model matrix X to the number of pedestrian casualties in collisions with buses, etc. The contingency table of injury counts between all mode pairings forms the “who hit whom” matrix for the city.

B.2 Making predictions

We use the same model equation to make predictions in hypothesised scenarios. The prediction equation requires us to specify the distances travelled in the scenario: call them \hat{M} and \hat{C} . Then we predict the expected number of injuries in the scenario, \hat{I} , as:

$$\hat{I} = \alpha \hat{M}^{\beta_1} \hat{C}^{\beta_2} \exp \left(\sum_{i=3}^P X_i \beta_i \right). \quad (38)$$

To aid interpretation, we can consider the ratio of the expected injuries in the scenario to the expected injuries in the baseline:

$$\frac{\hat{I}}{\mathbb{E}(I)} = \frac{\hat{M}^{\beta_1} \hat{C}^{\beta_2}}{M^{\beta_1} C^{\beta_2}} \quad (39)$$

$$= \left(\frac{\hat{M}}{M} \right)^{\beta_1} \left(\frac{\hat{C}}{C} \right)^{\beta_2}. \quad (40)$$

Then we can immediately read out, for example, that if M does not change ($\hat{M} = M$) then the fold change in injuries is equal to the fold change in cycling raised to the power β_2 : if $\beta_2 = 1$, then if cycling increases 25 times, so does the injury count. If $\beta_2 = 0.5$, then if cycling increases 25 times, the injury count increases five times.

B.3 β_1 and β_2 parameters

The question we need to answer is, given that we are using this model, what values should we choose for β_1 and β_2 ? This choice will impact on the other parameters to fit (α and β_i for $i \geq 3$) and, crucially, on the number of injuries we predict in scenarios.

Recall that there are multiple casualty modes and multiple “other party” modes, including NOV (no other vehicle). Another question we need to answer is how these values should differ for

different modes, in particular (a) where a mode's distance is not changing at all (or even very little) in scenarios, (b) where a mode is a combination of multiple modes, and (c) where there is no other mode.

C Power-law scaling relationships

The relationship described by Equation 2 closely resembles the family of power-law scaling models in e.g. Bettencourt et al. (2007). We could introduce this field of study to that family, describing safety in numbers as another relationship showing a power-law scaling relationship, meriting further investigation, for which solutions are occasionally proposed (Sim et al., 2015; Shalizi, 2011). This would be a very different way of approaching the problem, though there are some insights that might be useful in and of themselves.

Following the publication of Bettencourt et al. (2007) and other, similar work, some useful suggestions were made for studies such as these. Particularly relevant for our application is the requirement to fit models to per capita data (see Figure 19), if that is the quantity of interest to us (Shalizi, 2011). In the same work and in others (Leitão et al., 2016) we can find methods for investigating these types of data – perhaps there exists a statistical model already that would work for us. Also presented are various methods for robust specification and testing of null hypotheses, and discussion around competing models and in which circumstances one might conclude that “scaling is nonlinear” (Leitão et al., 2016).

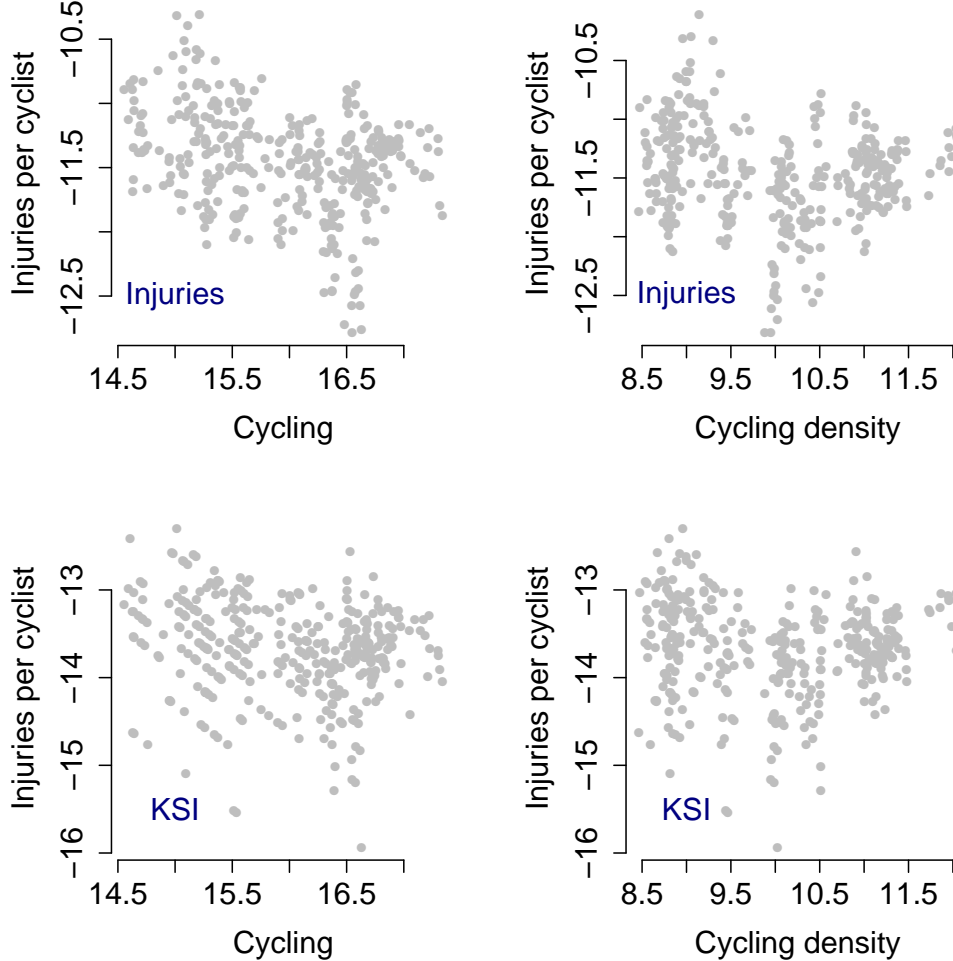


Figure 19: London injuries per cyclist distance, as a function of cyclist distance (left) and cyclist density (right). Cyclist density is calculated as total distance travelled divided by total distance of A, B and minor roads.

Here are some questions, specifically with that field in mind:

- (1) Is there a power-law scaling occurring in road injuries? Considering numbers alone (if we want

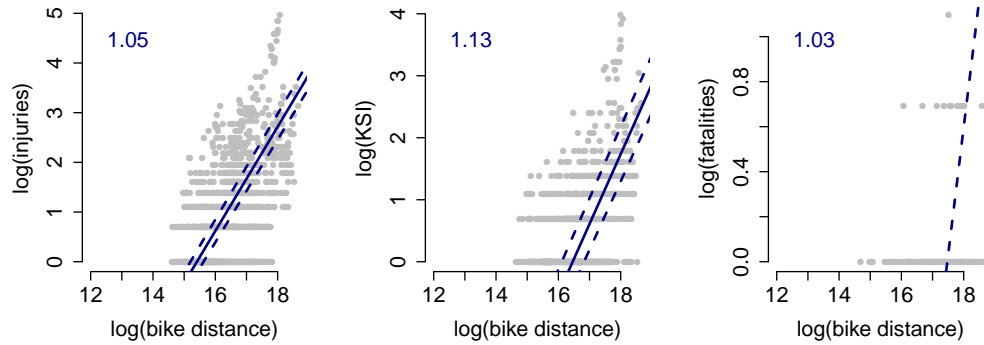


Figure 20: As in Figure 1, counting only injuries that occur in incidents involving no other vehicle (NOV).

to consider numbers at all), it seems to depend on exactly what we measure and exactly how we model it, e.g. injuries vs. KSI vs. fatalities (see e.g. Table 2). What does that mean for our interpretation? Study of no-other-vehicle casualties might be relevant here (Figure 1 vs. Figure 20).

- (2) Within a power-law scaling model, can we elicit the contributions of different modes? Could we apply a latent variable model, and learn the effects of fluctuations of mode types on the outcome?
- (3) Can we identify local (protective) effects?
- (4) If yes, can we design a mechanism that links local effects to the power-law scaling model, i.e. explains the deviance of $\beta_1 + \beta_2$ from 1 as a function of system size?

References

- Aldred, R., Goel, R., Woodcock, J. and Goodman, A. (2017), ‘Contextualising Safety in Numbers: a longitudinal investigation into change in cycling safety in Britain, 1991–2001 and 2001–2011’, *Injury Prevention* pp. injuryprev–2017–042498.
URL: <http://injuryprevention.bmj.com/lookup/doi/10.1136/injuryprev-2017-042498>
- Bettencourt, L. M. A., Lobo, J., Helbing, D., Kühnert, C. and West, G. B. (2007), ‘Growth, innovation, scaling, and the pace of life in cities’, *PNAS* **104**(17), 7301–7306.
- Birrell, P. J., Ketsetzis, G., Gay, N. J., Cooper, B. S., Presanis, A. M., Harris, R. J., Charlett, A., Zhang, X.-S., White, P. J., Pebody, R. G. and De Angelis, D. (2011), ‘Bayesian modeling to unmask and predict influenza A/H1N1pdm dynamics in London’, *Proceedings of the National Academy of Sciences* **108**(45), 18238–18243.
URL: <http://www.pnas.org/cgi/doi/10.1073/pnas.1103002108>
- Clauset, A., Shalizi, C. R. and Newman, M. E. (2009), ‘Power-law distributions in empirical data’, *SIAM Review* **51**(4), 661–703.
- de Sá, T. H., Tainio, M., Goodman, A., Edwards, P., Haines, A., Gouveia, N., Monteiro, C. and Woodcock, J. (2017), ‘Health impact modelling of different travel patterns on physical activity, air pollution and road injuries for São Paulo, Brazil’, *Environment International* **108**(June), 22–31.
URL: <http://dx.doi.org/10.1016/j.envint.2017.07.009>
- Elvik, R. and Bjørnskau, T. (2017), ‘Safety-in-numbers : A systematic review and meta-analysis of evidence’, *Safety Science* **92**(0349), 274–282.
URL: <http://dx.doi.org/10.1016/j.ssci.2015.07.017>
- Garder, P., Leden, L. and Pulkkinen, U. (1998), ‘Measuring the safety effect of raised bicycle crossings using a new research methodology’, *Transportation Research Record* **1636**(98), 64–70.
URL: <http://nacto.org/wp-content/uploads/2010/08/Measuring-the-Safety-Effect-of-Raised-Bicycle-Crossings-Using-a-New-Research-Methodology.pdf>
- Geyer, J., Raford, N., Pham, T. and Ragland, D. (2006), ‘Safety in numbers: data from Oakland, California’, *Transportation Research Record: Journal of the Transportation Research Board* **1982**, 150–154.
URL: <http://dx.doi.org/10.3141/1982-20>
- Kiers, H. A. and Smilde, A. K. (2007), ‘A comparison of various methods for multivariate regression with highly collinear variables’, *Statistical Methods and Applications* **16**(2), 193–228.
- Kleiber, M. (1947), ‘Body size and metabolic rate’, *Physiological Reviews* **27**(4).
- Leden, L. (2002), ‘Pedestrian risk decrease with pedestrian flow. A case study based on data from signalized intersections in Hamilton, Ontario’, *Accident Analysis and Prevention* **34**, 457–464.
- Leitão, J. C., Miotto, J. M., Gerlach, M. and Altmann, E. G. (2016), ‘Is this scaling nonlinear?’, *Royal Society open science* **3**, 150649.
URL: <http://dx.doi.org/10.1098/rsos.150649>
- Miranda-Moreno, L., Strauss, J. and Morency, P. (2011), ‘Disaggregate exposure measures and injury frequency models of cyclist safety at signalized intersections’, *Transportation Research Record: Journal of the Transportation Research Board* **2236**, 74–82.
URL: <http://trrjournalonline.trb.org/doi/10.3141/2236-09>
- Nordback, K., Marshall, W. E. and Janson, B. N. (2014), ‘Bicyclist safety performance functions for a U.S. city’, *Accident Analysis and Prevention* **65**, 114–122.
URL: <http://dx.doi.org/10.1016/j.aap.2013.12.016>

- Prato, C. G., Kaplan, S., Rasmussen, T. K., Hels, T., Giacomo, C., Kaplan, S. and Rasmussen, T. K. (2016), ‘Infrastructure and spatial effects on the frequency of cyclist-motorist collisions in the Copenhagen Region’, *Journal of Transportation Safety & Security* **8**(4), 346–360.
URL: <http://dx.doi.org/10.1080/19439962.2015.1055414>
- Schepers, J. P. and Heinen, E. (2013), ‘How does a modal shift from short car trips to cycling affect road safety?’, *Accident Analysis and Prevention* **50**, 1118–1127.
URL: <http://dx.doi.org/10.1016/j.aap.2012.09.004>
- Schepers, J. P., Kroeze, P. A., Sweers, W. and Wüst, J. C. (2011), ‘Road factors and bicycle – motor vehicle crashes at unsignalized priority intersections’, *Accident Analysis and Prevention* **43**(3), 853–861.
URL: <http://dx.doi.org/10.1016/j.aap.2010.11.005>
- Shalizi, C. R. (2011), ‘Scaling and hierarchy in urban economies’, *arXiv* p. 1102.4101v2.
- Sim, A., Yaliraki, S. N., Barahona, M. and Stumpf, M. P. H. (2015), ‘Great cities look small’, *Journal of The Royal Society Interface* **12**(109), 20150315.
URL: <http://rsif.royalsocietypublishing.org/lookup/doi/10.1098/rsif.2015.0315>
- Stumpf, M. P. H. and Porter, M. (2012), ‘Critical truths about power laws.’, *Science (New York, N.Y.)* **335**(6069), 665–6.
URL: <http://www.ncbi.nlm.nih.gov/pubmed/22323807%5Cnhttp://intersci.ss.uci.edu/wiki/pdf/Science-2012-Stumpf-665-6.pdf>