

Introducing spatial microsimulation with R: a practical

Lovelace, Robin
r.lovelace@leeds.ac.uk

April 24, 2014

Contents

1	Foundations	2
1.1	Prerequisites for this practical	2
1.2	Learning by doing	2
1.3	Some input data	3
1.4	The IPF equation	4
1.5	Test your understanding	5
1.6	IPF in a spreadsheet	6
1.7	IPF in R: a simple example	6
1.8	Iterations	8
2	CakeMap: A worked example	8
2.1	Loading and preparing the input data	8
2.2	Performing IPF on the CakeMap data	8
2.3	Processing the output	8
2.4	Mapping the results	8
3	Applying the technique in the real world	8
3.1	Model checking and validation	8
3.2	Visualising the results	9
3.3	Analysis	10
4	Discussion: learning outcomes, limitations and further work	11
5	Glossary	11
6	References	11

1 Foundations

1.1 Prerequisites for this practical

This practical is about spatial microsimulation in the software R. We suggest you install and take a look at this powerful program before getting started. We recommend using R within RStudio, which makes using R much easier. Instructions to install both R and RStudio can be found online: rstudio.com/ide/download/desktop.

The other prerequisite for the course is downloading the example data. These can be downloaded in a single zip file which can be found on the course’s GitHub repository: github.com/Robinlovelace/smsim-course. Click on the “Download ZIP” button to the right of this page and extract the folder into your desktop or other suitable place on your computer.

Once the folder has been successfully extracted open it in your browser and take a look around. You will subfolders entitled ‘cakeMap’, ‘data’, ‘figures’ and ‘rcode’. When you get to the sections that use R code, it is useful for R to operate from within the smsim-course-master folder. Probably the best way to set this up is to open the file ‘smsim-course.Rproj’ from within RStudio (fig. 1). Try this now and click on the ‘Files’ tab in the bottom right hand window of RStudio. Before using the power of R in RStudio it’s worth understanding a bit about ‘IPF’, the algorithm we will use to generate the synthetic population or ‘spatial microdata’ (fig. 2).

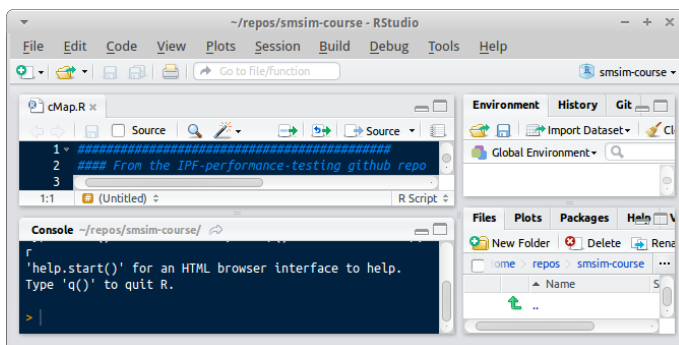


Figure 1: The RStudio user interface. Note the project title ‘smsim-course’ in the top right and the ‘Files’ tab at the top of the left hand window.

1.2 Learning by doing

As Kabacoff (2011, xxii) put it regarding R, “the best way to learn is to experiment” and the same applies to spatial microsimulation. We believe you will learn the technique/art best not by reading about it, but by doing it. Mistakes are inevitable in any challenging task and should not be discouraged. In fact it is by making blunders, identifying and then correcting them that many people learn best. Think of someone learning to skate: no one ever picks up a skateboard for the first time being able to ‘surf the sidewalks’. It takes time, patience and plenty of falls before you master the art. The same applies to spatial microsimulation.

One of the tricky things about spatial microsimulation for newcomers is its use of specialist language. It is important to know exactly what is meant by ‘special microdata’ and other technical terms. To this end we have created a glossary that provide succinct definitions (see section 5). Any term that is *italicised* in the text has a glossary entry.

Spatial microsimulation works by taking *microdata* at the individual level and using aggregate-level constraints to allocate these individuals to zones. The two main methods are *deterministic reweighting* and *combinatorial optimisation*. This practical takes the former approach using a process called *iterative proportional fitting* (IPF). IPF is used to increase the weights of individuals who are representative of the target area and reduce the weights of individuals who are relatively rare (Lovelace and Ballas, 2013). The output is a *spatial microdataset*.

A schematic of the process is shown in fig. 2. Take a look at the image and think about the process. But we don’t want to get bogged down in theory or applications in this course: we want to ‘get our hands dirty’. Next we will do just that, by applying the process to some example data.

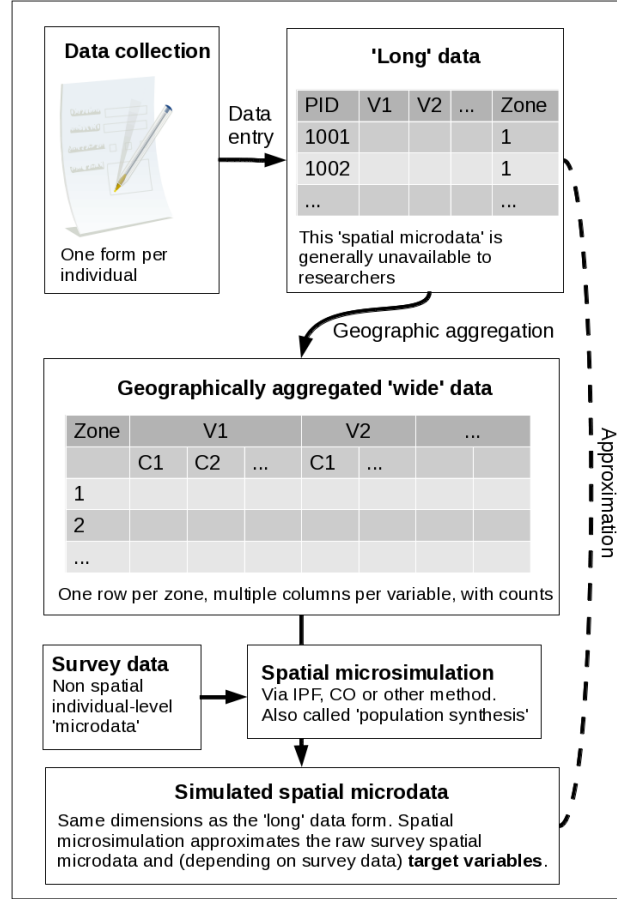


Figure 2: Schema of iterative proportional fitting (IPF) and combinatorial optimisation in the wider context of the availability of different data formats and spatial microsimulation.

1.3 Some input data

Let's start with a very basic dataset. To aid understanding, we will first do the reweighting by hand to understand the process before automating it on the computer. Table 1 shows 5 individuals, who are defined by two constraint variables: age and sex. Table 2 contains data for a hypothetical area. Table 3 illustrates this aggregate level table in a different form, to show our ignorance of interaction between age and sex. Finally, table 4 presents the hypothetical microdata in aggregated form, that can be compared directly to the aggregate data presented in Table 3.

Table 1: A hypothetical input microdata set

Individual	Age	Sex	Weight
1	59	Male	1
2	54	Male	1
3	35	Male	1
4	73	Female	1
5	49	Female	1

Table 2: Small area constraints (s).

Constraint \Rightarrow	i		j	
Category \Rightarrow	i_1	i_2	j_1	j_2
Area \Downarrow	Under-50	Over-50	Male	Female
1	8	4	6	6

Table 3: Small area constraints expressed as marginal totals, and the cell values to be estimated.

	Marginal totals	j		
	Age/sex	Male	Female	T
i	Under-50	?	?	8
	Over-50	?	?	4
	T	6	6	12

Table 4: The aggregated results of the weighted microdata set ($m(1)$). Note, these values depend on the weights allocated in Table 1 and therefore change after each iteration

	Marginal totals	j		
	Age/sex	Male	Female	T
i	Under-50	1	1	2
	Over-50	2	1	3
	T	3	2	5

1.4 The IPF equation

Using these tables we readjust the weights of the individuals so that their sum equals the total population of the area presented in Table 3. For each constraint, the weights are multiplied by the aggregate level values from Table 2 and divided by the respective marginal total of the microdata (see table 4). This is done one constraint at a time, as described in fig. 3 and more formally by eq. (1) for constraint i (age in this case):

The new weight of an individual ($w(n+1)_{ij}$) is equal to
 Their current weight ($w(n)_{ij}$) multiplied by
 The frequency of their category
 After applying the previous constraint (sT_i)
 Divided by the value of their category
 In the census/constraint table ($mT(n)_i$)

Figure 3: Equation (1) in word form.

$$w(n+1)_{ij} = \frac{w(n)_{ij} \times sT_i}{mT(n)_i} \quad (1)$$

where $w(n+1)_{ij}$ is the new weight for individuals with characteristics i (age, in this case), and j (sex), $w(n)_{ij}$ is the original weight for individuals with these characteristics, sT_i is element marginal total of the small area constraint, s (Table 2) and $mT(n)_i$ is the marginal total of category j of the aggregated results of the weighted microdata, m (Table 4). n represents the iteration number.

Do not worry about understanding the above equation for now. More important is implementing it. Follow the emboldened values in tables 1 to 4 to see how the new weight of individual 3 is calculated for the age constraint. Table 5 illustrates the outcome. Notice that the sum of the weights is equal to the total population, from the constraint variables.

Table 5: Reweighting the hypothetical microdataset in order to fit Table 2.

Individual	Sex	age-group	Weight	New weight, $w(2)$
1	Male	Over-50	1	$1 \times 4/3 = \frac{4}{3}$
2	Male	Over-50	1	$1 \times 4/3 = \frac{4}{3}$
3	Male	Under-50	1	$1 \times 8/2 = 4$
4	Female	Over-50	1	$1 \times 4/3 = \frac{4}{3}$
5	Female	Under-50	1	$1 \times 8/2 = 4$
Total			5	12

After the individual level data have been re-aggregated (table 6), the next stage is to repeat eq. (1) for the age constraint to generate a third set of weights, by replacing the i in sT_i and $mT(n)_i$ with j and incrementing the value of n :

$$w(3)_{ij} = \frac{w(2)_{ij} \times sT_j}{mT(2)_j} \quad (2)$$

1.5 Test your understanding

To test your understanding of IPF, try to apply eq. (2) to the information above and that presented in table 6. This should result in the following vector of new weights, for individuals 1 to 5. Calculate the correct values and pencil them in in place of the question marks.

$$w(3) = \left(\frac{6 \times \frac{4}{3}}{6\frac{2}{3}}, \frac{6 \times \frac{4}{3}}{?}, \frac{4 \times 6}{6\frac{2}{3}}, \frac{?}{?}, \frac{4 \times 6}{5\frac{1}{3}} \right) \quad (3)$$

After simplifying these fractions, the results are as follows. One ‘sanity’ check of your method here is whether the sum of these weights is still equal to the area’s total population of twelve. Test this is true:

$$w(3) = \left(\frac{6}{5}, \frac{6}{5}, \frac{18}{5}, \frac{3}{2}, \frac{9}{2} \right) \quad (4)$$

What do the weights in $w(3)$ actually mean? They indicate how representative each individual is of the target zone after one iteration of IPF, constraining by age and sex. Individual number 5 has the highest weight because there is only one young female in the survey dataset yet seemingly many in the area in question.

Notice also that after each iteration the fit between the marginal totals of m and s improves.¹

Table 6: The aggregated results of the weighted microdata set after constraining for age ($m(2)$).

Marginal totals	Age/sex	i		
		Male	Female	T
j	Under-50	4	4	8
	Over-50	$\frac{8}{3}$	$\frac{4}{3}$	4
	T	$6\frac{2}{3}$	$5\frac{1}{3}$	12

¹This can be checked by comparing the aggregated weighted individuals with the small area constraints. Total absolute error (TAE), defined as the sum of all differences between simulated and observed marginal totals, improves between $m(1)$ to $m(2)$, falling from 14 to 6 in table 4 and table 6 above. TAE for $m(3)$ (not shown, but calculated by aggregating $w(3)$) improves even more, to 1.3. This number would eventually converge to 0 through subsequent iterations, a defining feature of IPF.

1.6 IPF in a spreadsheet

To ensure smooth transition between the IPF process described in mathematics above and its implementation in R this section provides an intermediary stage: the speed of computation and the visual support of a graphical user interface (GUI).

1.7 IPF in R: a simple example

So far we have implemented IPF by hand and in a spreadsheet. This section explains how the IPF *algorithm* described above is implemented in R, using the same input data.² The code snippets are taken from the ‘simple.R’ script file (just a plain text document, like all R scripts) so you can copy and paste. However, it is highly recommended to not look at this file until you need to, towards the end of section. This is ‘leaning by typing’!

The data presented in the above worked example are saved in the ‘simple’ sub-folder of ‘data’ as .csv files. To load them in R, we use the following commands:

Listing 1: Loading the input data in R

```
ind <- read.csv("data/simple/ind.csv") # load the individual level data
cons <- read.csv("data/simple/cons.csv") # load aggregate constraints
```

To check that the *data frames* have loaded correctly, simply ask R to print their contents by entering the object’s name. To display the individual level data, for example simply type the following:³

Listing 2: Checking the contents of the individual level data frame

```
ind # Show the data frame in R
##  id age sex
##   1  59   m
##   2  54   m
##   3  35   m
##   4  73   f
##   5  49   f
```

Note that the input data is identical to the tables illustrated above for consistency. Instead of constraining for just one zone as we did in the example above, we will fit the data to six zones here. To subset a dataset, add square brackets to the end of the object’s name in R. Entering `cons[1:2,]`, for example should output the first two rows of the constraint variable: within the square brackets the number before the comma refers to columns (blank means all columns) and the numbers after the comma refer to rows. If the output of this command looks like the text below, congratulations are in order: you have successfully loaded the constraint dataset.

Listing 3: Printing a subset of the constraint data in R

```
cons[1:2,]
##      X16.49 X50.  m f
##   1         8    4 6 6
##   2         2    8 4 6
```

²A fuller tutorial is available from Rpubs, a site dedicated to publishing R analyses that are reproducible. This tutorial uses the RMarkdown mark-up language, which enables R code to be run and presented within documents. See <http://rpubs.com/RobinLovelace/5089>.

³An alternative way to do this is to click on the object’s name in the ‘Environment’ tab in RStudio’s top right window.

Note that the top row is identical to table 2 — we can therefore compare the results of doing IPF on the computer with the results obtained by hand.

To make the individual level data comparable to the aggregate level we must first categorise it and count the number in each category. This job is performed by the R *script* `categorise.R`, contained in the ‘simple’ data folder. Run it by entering `source("data/simple/categorise.R")`. This results in the creation of a new object called `ind.cat`, which converts the variables into a matrix of 0s and 1s with same number of columns as the constraint data. Note that the ‘`categorise.R`’ script must be modified each time it is used with a different dataset. To check that the script has worked properly, let's count the number of individuals it contains:

Listing 4: Output from `ind.cat`

```
colSums(ind.cat)
##      a16.49      a50.      m      f
##           2         3         3         2
```

The sum of both age and sex variables is 5 (the total number of individuals): it worked! Now the data is in the same form as the constraint variables and we have checked the data makes sense, we can create the ‘weight’ matrix and begin reweighting the individuals to zones:⁴

Listing 5: Creating the weight and aggregated individual level matrices

```
weights <- array(1, dim=c(nrow(ind),nrow(cons)))
ind.agg <- matrix(rep(colSums(ind.cat), nrow(cons)), nrow(cons))
```

With all of these objects in place, we are ready to begin allocating new weights to the individuals via IPF. Below is written in code the IPF formula presented in section 1.4.

Listing 6: Creating new weights for the individuals based on constraint 1 (age) via IPF

```
for (j in 1:nrow(cons)){
  for(i in 1:ncol(con1)){
    weights[which(ind.cat[,i] == 1),j] <- con1[j,i] / ind.agg[j,i]}}
```

The above code creates the weight matrix that allocates individuals to zones. Aggregating this data to zones allows us to update ‘`ind.agg`’ with new values that will be closer to the constraint variables. This is verified in the ‘`simple.R`’ script file with a line of code that calculates the Total Absolute Error after each iteration — `sum(sqrt((ind.agg-cons)^2))`. Check the error decreases with each iteration.

At this stage you should switch to running the code in ‘`simple.R`’. In RStudio to run a line of code from the top left window, simply press **Ctrl-Enter**. Use this technique to constrain by the second constraint, sex. Notice that we have just done in code what was previously done by hand. Thus we can verify the output of the R implementation against that obtained through mental arithmetic. Compare the following listing with eq. (4): are the weights the same?

⁴The above code will probably seem daunting and not make sense unless you are an experienced R user. Do not worry about this for now, you can always check what is going on later by checking R’s documentation (e.g. try entering ‘`?rep`’) or by pressing the tab key when entering arguments for R commands. For the time being, it is best to press-on with the example and understand the concepts — the details can be looked at later.

Listing 7: $W(3)$ after constraining by age and sex in R. Compare with eq. (4)

```
weights3[,1] # check the weights allocated for zone 1
##      [1] 1.2 1.2 3.6 1.5 4.5
```

1.8 Iterations

In the above example, we have seen the ‘bare bones’ of spatial microsimulation using IPF to generate weights from which sample populations can be created. To perform multiple iterations of the same model, we have prepared a slightly more complex script than ‘simple.R’ called ‘etsim.R’ that can be found in the ‘data/simple/’ folder. The etsim script generates weights for as many iterations as you please (just change the value of `num.its`). This happens by calling another script called ‘e2.R’ which is simply run repeatedly for each iteration.

There is great scope for taking the analysis further: some further tests and plots are presented on the on-line versions of this section. The simplest case is contained in Rpubs document rpubs.com/RobinLovelace/6193 and a more complex case (with three constraints) can be found in Rpubs document 5089. For now, however, we progress to a more complex example, CakeMap.

2 CakeMap: A worked example

In this section we will process real data to arrive at an important result: an estimate of the amount of cake eaten in different wards in Leeds. The example is deliberately rather absurd: hopefully this will make the technique more memorable. An attempt has been made to present the method in a generalisable way as possible, allowing users to apply the technique described here to their own data.

The code needed to run the main part of the example is contained within ‘cakeMap.R’. Note that this script makes frequent reference to files contained in the ‘cakeMap’ folder ‘data’. `read.csv("data/cakeMap/ind.csv")`, for example, is used to read the individual level data into R.

Because this example uses real world data (albeit in anonymised form), we will see how important the process of ‘data cleaning’ is. Rarely are datasets downloaded in exactly the right form for them to be pumped directly into a spatial microsimulation model. For that reason we describe the process of loading and preparing the input data first: similar steps will be needed to use your own data as an input into a spatial microsimulation model.

2.1 Loading and preparing the input data

2.2 Performing IPF on the CakeMap data

2.3 Processing the output

2.4 Mapping the results

3 Applying the technique in the real world

Here we will continue with the cakeMap model we ran in the previous section.

3.1 Model checking and validation

To make an analogy with food safety standards, openness about mistakes is conducive to high standards (Powell et al., 2011). Transparency in model verification is desirable for similar reasons. The two main strategies are 1) comparing the model results with knowledge of how it *should* perform *a-priori* (model checking) and 2) comparison between the model results and empirical data (validation).

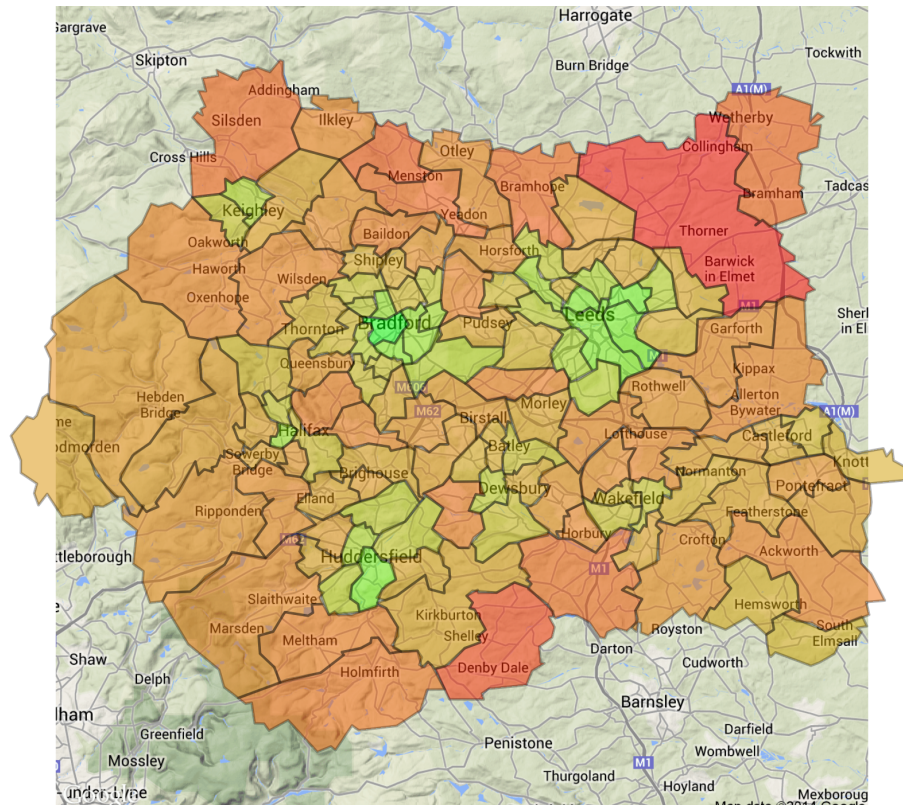


Figure 4: Choropleth map of the spatial distribution of average frequency of cake consumption in West Yorkshire, based on simulated data.

A proven method of checking that data analysis and processing is working is wide ranging and continual visual exploration of its output (Janert, 2010).

Beyond typos or simple conceptual errors in model code, more fundamental questions should be asked of spatial microsimulation models. The validity of the assumptions on which they are built, and the confidence one should have in the results are important.

3.2 Visualising the results

Visualisation is an important part of communicating quantitative data, especially so when the datasets are large and complex so not conducive to description with tables or words.

Because we have generated spatial data, it is useful to create a map of the results, to see how it varies from place to place. The code used to do this found in ‘cMapPlot.R’. A vital function within this script is ‘merge’, which is used to add the simulated cake data to the geographic dataframe:⁵

Listing 8: The merge function for joining the spatial microsimulation results with geographic data. Compare with eq. (4)

```
merge(wardsF, wards@data, by = "id")
```

The above line of code by default selects all the data contained in the first object (‘wardsF’) and adds to it new variables from the second object based on the linking variable (in this case “id”). Also in that script file you will encounter the function **fortify**, the purpose of which is to convert the spatial data object into a data frame. More on this process is described in ?. The final map result of ‘cakeMapPlot.R’ is illustrated below (fig. 4).

⁵‘join’ is an alternative to merge from the ‘plyr’ package also used in the ‘cMapPlot.R’ script that performs the same task. Assuming ‘plyr’ is loaded — ‘library(plyr)’ you can read more about join by entering ‘?join’ in R.

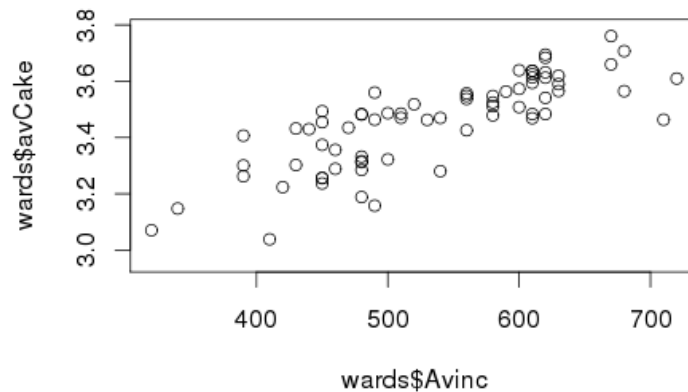


Figure 5: Relationship between modelled average ward income and simulated number of cakes eaten per person per week.

3.3 Analysis

Once a spatial microdataset has been generated that we are happy with, we will probably want to analyse it further. This means exploring it — its main features, variability and links with other datasets. To illustrate this process we will load an additional dataset and compare it with the estimates of cake consumption per person generated in the previous section at the ward level.

The hypothesis we would like to test is that cake consumption is linked to deprivation: More deprived people will eat unhealthily and cake is a relatively cheap ‘comfort food’. Assuming our simulated data is correct — a questionable assumption but let’s roll with it for now — we can explore this at the ward level thanks to a dataset on modelled income from neighbourhood statistics.

Because the income dataset was produced for old ward boundaries (they were slightly modified for the 2011 census), we cannot merge with the spatial dataset based on the new zone codes. Instead we rely on the name of the wards. The code below provides a snapshot of these names and demonstrates how they can be joined using the ‘join’ function.

Listing 9: The merge function for joining the spatial microsimulation results with geographic data. Compare with eq. (4)

```
wards@data <- join(wards@data, imd)
summary(imd$NAME %in% wards$NAME)
##      Mode  FALSE  TRUE  NA's
## logical    55    71     0
```

The above code first joins the two datasets together and then checks the result by seeing how many matches names there are. In practice the fit between old names and new names is quite poor: only 71 out of 124. In a proper analysis we would have to solve this problem (e.g. via the command `pmatch` which stands for partial match). For the purposes of this exercise we will simply plot income against simulated cake consumption to gain a feeling what it tells us about the relationship between cake consumption and wealth (fig. 5).

4 Discussion: learning outcomes, limitations and further work

What have we learned during the course of this practical tutorial? Given that it was geared towards beginners, from one perspective one could say ‘not a lot’ relative to the many potential applications. The emphasis on understanding of the basics was deliberate: it is only on strong foundations that large and long-lived buildings can be constructed and the same applies to geographical research.

The focus on strong foundations is also sensible because the number of people actively involved in spatial microsimulation research, let alone development of new software and methods, is very small. Specialist skills and understanding is needed to join this research area and a central outcome of this course has been to provide these. More specifically, we have learned how to:

- Perform iterative proportional fitting, one of the commonly used techniques for spatial microsimulation, by hand, in a spreadsheet and in R.
-
-

5 Glossary

- **Algorithm:** a series of computer commands which are executed in a well defined order. Algorithms process input data and produce an output.
- **Data frame:** a type of object (formally referred to as a class) in R, data frames are square tables composed of rows and columns of information. As with many things in R, the best way to understand dataframes is to create them and experiment. The following creates a dataframe with two variables: name and height:

```
data.frame(name = c("Robin", "Phil"), height.cm = c(172, 174))
```

Note that each new variable is entered using the command `c()`. This is how R creates objects with the *vector* data class — a one dimensional matrix — and that text data must be entered in quote marks.

- **Iteration:** one instance of a process that is repeated many times until a predefined end point, often within an *algorithm*.
- **Iterative proportional fitting (IPF):** an iterative process implemented in mathematics and algorithms to find the maximum likelihood of cells that are constrained by multiple sets of marginal totals. To make this abstract definition even more confusing, there are multiple terms which refer to the process, including ‘biproportional fitting’ and ‘matrix raking’. In plain English, IPF in the context of spatial microsimulation can be defined as *a statistical technique for allocating weights to individuals depending on how representative they are of different zones*. IPF is a type of deterministic reweighting, meaning that random numbers are not needed to generate the result and that the output weights are real (not integer) numbers.
- **Iteration**

6 References

- Janert, P.K., 2010. Data analysis with open source tools. O’Reilly Media.
- Kabacoff, R., 2011. R in Action. Manning Publications Co.

- Lovelace, R., Ballas, D., 2013. Truncate, replicate, sample: A method for creating integer weights for spatial microsimulation. *Computers, Environment and Urban Systems* 41, 1–11.
- Powell, D.a., Jacob, C.J., Chapman, B.J., 2011. Enhancing food safety culture to reduce rates of foodborne illness. *Food Control* 22, 817–822.