

# Econ 21410 - Problem Set II

## Schelling's segregation and related models\*

April 30, 2015

This homework should be done in LaTeX. The homework will be graded on correctness, but will also heavily weight clarity and professionalism. Being able to produce clean, clear, well documented write-ups and code is an important skill which will be rewarded. Its better to not do some of the harder parts than to turn in an incomprehensible document. Your R script as well as a log-file should be submitted. Alternatively, use knitr to print your code inline in your latex document.

Make sure to write code which is clear and flexible. Read the whole problem before you begin coding. Some parameters will change and the code should be written in a way to make this easy to implement. We will re-use code in this course. Flexibility and documentation now will save you headaches later in the quarter. Remember to properly indent your code!

**SUBMISSION:** The homework must be emailed to Oliver and myself by 9:30a.m. Monday, April the 13th. The email must include a pdf with the filename `lastname_pset2.pdf` and R code called `lastname_pset2_code.R` where "lastname" should be replaced with your last name. The subject of your email should be [ECON 21410: pset2 submission]

Remember that asking and answering questions on our github page, coming to office hours to ask questions, and contributing to the class wiki are all worth participation credit, which is 10% of your grade in this class.

## 1 Writing functions and data types in R

Sometimes it is helpful to write a set of functions that are generally useful for data analysis that you will use across projects. This is also a great way to increase your understanding of R if you are already a good programmer.

Write a function that takes a specific variable from a data.frame and does the following: (note: this problem is pretty technical and will not be weighted as heavily in the grading. Please focus on the Schelling section below first.)

1. If the variable is of type "character" (i.e., a string) return 5 random examples. Return the proportion that have the value "" (an empty string).
2. If the variable is of type "factor" return the levels of the factor and the proportion of each data that is at each level

---

\*Please email [johneric@uchicago.edu](mailto:johneric@uchicago.edu) and [obrowne@uchicago.edu](mailto:obrowne@uchicago.edu) if you have questions.

3. If the variable is of type “numeric” return the mean, median, mode, min, max, 25th percentile, and 75th percentile.
4. Regardless of type, return the proportion of the variables that are of value NA
5. Regardless of the type, return the number of non-missing values
6. (optional) If passed a numeric or character type, check the number of unique values and if its less than 100, convert the variable into a factor inside your code and output results for factor plus a string saying its actual type.
7. (optional) Make the 100 above an option set by the function
8. (optional) Rewrite the function above so that it takes the name of a data.frame and the name of a variable in the data.frame as separate arguments
9. (optional) Rewrite the function above using data.tables and the more efficient data.table manipulations
10. (optional) Rewrite the function above to take an optional “by” argument that will change the function to output the same results, but conditional on every value in the ”by” argument. An example would be to pass the value ”year” to the function, and have all the results passed by year.
11. (optional) Make the function return an S3 object with a print() method that displays the output of the function in a pleasant way
12. ( I will award 0.25 side-project points for each completed optional part above. Please make sure to note any side projects you complete at the end of the pset. )

## 2 Schelling’s Segregation Model

Consider a version of Schelling’s Segregation Model implemented in the following way:

1.  $\{n_1, \dots, n_N\}$  are generated where  $n^r$  are red,  $n^g$  are green, and  $n^b$  are blue.
2. Each individual is initially placed (uniformly) at random on  $[0, 1] \times [0, 1]$
3. Start with individual  $n_1$  (who has type  $t \in (r, b, g)$ ) then proceed with the algorithm.
  - a) Take individual  $n_{i,t}$ . If *at least*  $j_t$  of their  $m_t$  closest neighbors are the same color as them, move on to  $n_{i+1}$ ,
  - b) if fewer than  $j_t$  of their  $m_t$  closest neighbors are the same color as them, randomly draw a new living location and move on to  $n_{i+1}$ .
  - c) Continue this process until no individuals remain who wish to move.

## Machinery for the Schelling Model:

1. write a function that calculates the distances between a coordinate point  $(x_i, y_i)$  and a vector of coordinate points  $[(x, y)]$ . This should return a vector of distances
2. Write a function which simulates Schelling's Segregation model. The function should take the size of each population,  $j_t$  and  $m_t$  for each type. The function should allow for up to three types.
  - The function should output a plot (or the data necessary to output a plot) of the initial distribution of agents and the final distribution of agents (even better if it outputs some intermediate plots)
  - The function should return the data for the final allocation of individuals.
  - The function should return the number "cycles" the algorithm takes
3. We would like to study the "amount" of segregation in an our "city". To do this, we will write a function which will compute three different segregation metrics. The function should take the simulated final data from your Schelling function and return the three following metrics (focus on the first metric, questions involving the second two indexes will in total be not be worth more than 5% of the grade and are more challenging):
  - Similar neighbors index: For each individual, calculate the proportion of their  $m_t$  nearest neighbors that are the same type as them. Take the average of this number across individuals.
  - Dissimilarity index: grid up the  $[0, 1] \times [0, 1]$  city into "blocks" of size  $0.2 \times 0.2$ . Using these blocks have the function return the dissimilarity index.<sup>1</sup>
  - Gini index: Using the same grid as above, have the function return the gini index.

## Output for Schelling

1. Run a "baseline" model with two populations of 250. Let each group care about their 8 nearest neighbors and lets assume members from both groups are "happy" if half of their nearest 8 neighbors are the same color as them.
  - Make a plot of the initial distribution of individuals and the final distribution of individuals.
  - Run your model from a few different starting seeds and see how stable your results are above and discuss (2-5 sentences)
2. Make a plot showing how the number of iterations changes as you increase the populations of the two groups (symmetrically).

If you wanted to make this plot smoother you could consider plotting averages over several simulations

	pop.sizes	iter.by.popn	time.by.popn
1	50.00	16.00	1.46
2	100.00	12.00	1.79
3	150.00	16.00	3.92
4	200.00	16.00	5.18
5	250.00	15.00	5.82
6	300.00	16.00	8.03
7	350.00	21.00	12.50
8	400.00	21.00	14.59
9	450.00	15.00	12.14
10	500.00	19.00	17.63

Table 1: Convergence Time by Population

3. Make a table showing how run-time increases as you increase the populations of the two groups (symmetrically) (hint, see the command “system.time()”).
4. Calculate the similar-neighbor index, discuss (1-4 sentences).
5. Make a plot of how the similar-neighbor index changes as you increase the ratio of nearest neighbors that need to be of the same type for the individual to be happy from .1 to .9
6. Make a plot showing how the similar-neighbor index changes as you increase the number of individuals in each population from 50 to 500.
7. Make a plot showing how the similar-neighbor index changes based on the number of nearest neighbors considered (from 5 to 30).

**Differences between populations:** Now lets consider how the model changes when the two populations do not have the same characteristics.

1. Keeping the populations at 250, let the first population be happy if  $\frac{6}{8}$  of it’s nearest neighbors are of the similar type while the second population is happy if  $\frac{3}{8}$  of its nearest neighbors are of similar type.
2. Produce a plot showing the new allocation of individuals.
3. What is the value of the similar-neighbor index? (run the function a few times with different random seeds to see how much this varies). Briefly discuss.

We seem to observe higher levels of segregation here, which suggest that the aggregate level of segregation may be driven by the group with the stronger preferences to be nearby similar types.

4. Now let there be 500 of the first population, but only 100 of the second population. Again make a plot showing the final spatial allocation and calculate the similar-neighborhood index. How did this change? Discuss in 2-8 sentences.

When the sub-population with stronger preferences is more numeros, inequality measures seem to increase further

---

<sup>1</sup>Calculating the dissimilarity index and the gini index are a fair amount of additional work. These will be worth substantially fewer points than the similar neighbor index. Make sure you finish the rest of the problem set first.

**Extending the model to three populations** Now let's evaluate if the model changes when we introduce a third population.

1. Let each population have 150 individuals. Assume they are happy if  $\frac{3}{8}$  of their neighbors are the same as them.
2. Produce a plot showing the new final allocation of individuals.
3. What is the similar-neighborhood index? (run the function a few times with different random seeds to see how much this varies). Briefly discuss.

Even with a third population, we see measure similar levels of segregation as before. This is perhaps unsurprising since we observe very similar clustering patterns. On the other hand perhaps we would expect to see a lower Similar Neighbor Index because there are fewer individuals of each type.

4. Now let there be 500 of the first population, but only 100 in the second and third population. Let the first population be happy if  $\frac{9}{12}$  of its nearest neighbors are of the similar type while the second and third population are happy if  $\frac{2}{12}$  of its nearest neighbors are the same. What sort of model or scenario would this correspond to? Run this model several times and look at the final distribution of results. What is the "take away" of this model which we could use to test or inform our work with real data? Discuss (no more than a page).

We measure slightly lower levels of segregation here, but this is perhaps only because the two minorities are willing to intermingle. If we aggregated the two minorities into one I suspect you would measure similar levels of inequality as we did in the *Differences between populations* section. Note that this should tell us that we need to think carefully about how we aggregate subgroups.

You could analogize this simulation to a real world city with say *white, black and latino* populations. However before you did this you would have to stop and think very carefully about what would be the appropriate way to model the preferences of each group. Are real world cities consistent with the results of this simulation?

It is also worth thinking about how the results of this simulation compare to the Schelling grid simulations we did in class. In this model individuals can live arbitrarily and you will notice that these simulations lead to a lot of heterogeneity in density within cities. It seems that the majority in these simulations tend to live on average in higher densities in this simulation in order to avoid the minorities. Is this realistic? How would you measure the density of each subpopulation in this simulation?

**Alternative Segregation Indexes** Take your functions which calculate the dissimilarity-index and gini-index and compare how these vary compared to the similar-neighborhood index.<sup>2</sup>

### 3 Code Review

A surprisingly large part of coding is learning from and incorporating code others have already written. In this problem you will download code for a function which "simulates a peer-effects model". It's your job to back out how the model works and answer some general questions about the code.

---

<sup>2</sup>This is purposefully unstructured (like most real work). Figure out what is interesting about the differences and discuss.

1. Study the code and give an over-view of how this peer-effects model works. Don't discuss how the code works, but rather, sketch the model I used when implementing this code.

- Each individual is ordered from  $i = 1, \dots, 200$ , she is uniformly randomly placed at a point  $x_i, y_i$  in cartesian space on the unit interval  $U[0, 1] \times U[0, 1]$ . This individual also receives a preference shock  $\epsilon_i \sim U[-0.5, 0.5]$ .
- In the initial plot, each individual chooses to be red if their preference shock  $\epsilon > 0$
- In the middle plot individuals act sequentially from  $1, \dots, 200$ .  
Of all the individuals who have acted before individual  $i$  (i.e  $\forall j < i$ ). Individual  $i$  looks at the actions of the twenty nearest to him and calculates the fraction who chose to be red  $p_{red}$ <sup>3</sup>. Then individual  $i$  chooses to be red if

$$-\frac{1}{2.5} + \frac{2}{2.5}p_{red} + \epsilon_i > 0$$

- The second plot again iterates sequentially through all of the individuals. However now all individuals consider their 20 closest neighbors (regardless of where their order. This time they make a decision according to a slightly different criteria.

$$-\frac{1}{s} + \frac{2}{s}p_{red} + \epsilon_i > 0$$

2. Are some people more affected by peers than others in this model?

Yes. In two manners

- Firstly: In the middle plot, the first individual  $i = 1$  is not affected by his peers, and subsequent individuals  $i \in 2, \dots, 20$  will be influenced by a smaller number of peers than individuals  $i \in 21, \dots, 200$  individuals. This will make the individuals who act first relatively more important.
- Secondly: Individuals whose preference shocks  $\epsilon_i$  are larger in absolute value will be less influenced by their peers than individuals with  $\epsilon_i$  closer to zero

3. What do the for-statement on line 42 and the while-statement on line 57 do (1-2 sentences each)?

- In the for-loop on line 42, we move through all the individuals once and they each make a decision to choose red or not.
- In the while loop on line 57 is redundant. Since at the end of the loop we set `data.old = data` the loop will execute once and then be trivially satisfied

4. Why do I output three different plots in the code? What do they each show? Producing all three plots allows us to compare the three different models, with no peer effects, after iterating through with peer effects once, and after iterating through with peer effects a second time.

5. What do  $k$ ,  $n$ , and  $s$  do in the code?

- $k$ : The number of nearby individuals you are influenced by

---

<sup>3</sup>if  $i = 1$  then  $p_{red} = 0$ , if  $i < 20$  then individual  $i$  looks at everybody who came before him regardless of distance

- $n$ : The total number of individuals
  - $s \in (0, \infty)$ : Is a sensitivity parameter which determines the size of the peer effects, the larger  $s$  the smaller the peer effects are relative to the individual shocks
6. How does changing  $s$  and  $k$  affect the code (1-3 sentences)?  
 Changing  $s$  will decrease individual's sensitivity to the peer effects, this will result in a more random (less segregated) distribution of types.  
 Changing  $k$  will increase the number of people each individual is influenced by. It is not clear what impact this will have on segregation ex ante. Very large  $k$ 's will lead to large uniform areas of influence and so little segregation, but very small  $k$ 's will lead to small areas of influence and so this may not lead to much clustering.
  7. Reuse your similar-neighbor index function from above to calculate the degree of segregation in the baseline model I run in my .R file. How does this model's level of segregation compare to our baseline Schelling model?

## Research

- Suggest one research idea based around the models we considered in this homework (no more than 3 sentences).
- List three topics you may be interested in doing research on. These can be broad, such as "The returns to community college", or very narrow, such as "Advertising for the Xbox One and the PlayStation 4". As undergraduates, it's helpful to stick to things that really interest you or things you really care about (anything going on in your home state?).<sup>4</sup>
- Propose some research idea related to one of those topics (3-10 sentences).
- If you are a 4th year who wrote a BA or you may want to extend work you have already done for the final, please write a 3 sentence summary of your BA, then write 1-5 sentences on how you may extend it for this course.<sup>5</sup> You will need to get individual approval from Oliver or I if you wish to extend work you have already completed for the course final.

## Side Projects

1. Download census tract data on race for the city of Chicago. Calculate the gini and dissimilarity indexes using your function. If possible make a map in R of these results. Write up no more than 1 page (not including figures) discussing your results (3 points)
2. Rewrite some or all of the code for this problem in Julia or C++ with Rcpp<sup>6</sup> (up to 3 points)

---

<sup>4</sup>For example, as an undergraduate I could have maybe have tried to write a hedonic pricing model for resell of high-end acoustic guitars (maybe scraping the data off of sales sites), or I could have written a paper on how Alaska's economy is counter-cyclical with rest of the Nation and the resulting impact of United States' monetary policy on Alaska's economy.

<sup>5</sup>Broad ideas are fine. If you don't have a good answer, come talk to us at office hours sometime in the next week or two

<sup>6</sup>If you tackle Rcpp send me an email and I can try to help you get started. You will most likely want to use RcppArmadillo, which links to the fantastic Armadillo C++ library which has syntax somewhat similar to matlab

3. write a wiki entry with the examples for any of the following R commands: “str”, “%in%”, “match”, “head”, “subset”, “with”, “get”, “all.equal”, “identical”, “complete.cases”, “pmax”, “rownames”, “diag”, “sweep”, “rep”, “rep\_len”, “unlist”, “rev”, “choose”, “expand.grid”, “replicate”, “duplicated”, “unique”, “table”, “ftable”, “sort”, “order”, “rank”, “crossprod”, “eigen”, “solve”, “rcond”, “cat”, “message”, “dir” (0.5 points), no more than 1 per person per week, only one writeup per command, max 4 for any individual).
4. Improve the wiki entry on ggplot2 (and the wrapper qplot). This is an extremely powerful plotting tool. Please provide examples and their corresponding plots. (up to 1 point, but multiple people can contribute and extend with different examples, different plot types, etc)
5. Build a shiny app (<http://www.rstudio.com/shiny/>) which lets users configure and run the Schelling model (or the grid Schelling Model) using a graphical user interface (up to 2.5 points).
6. Use knitr to add animated plots to your pdfs (animations only work in newer versions of Adobe products) (1.5 points)
7. Read the Thomas Schelling’s 1971 and/or 1969 paper and write a short (2-5 page) summary and review of his paper(s). If you read both papers, discuss the difference between the two (up to 2.5 points). (<http://www.tandfonline.com/doi/pdf/10.1080/0022250X.1971.9989794>) (<http://www.jstor.org/stable/pdfplus/1823701.pdf>)

## Useful Resources

[https://www.census.gov/hhes/www/housing/resseg/pdf/app\\_b.pdf](https://www.census.gov/hhes/www/housing/resseg/pdf/app_b.pdf)