

Open research methods in computational social sciences and humanities: introducing R

Markus Kainu

January 23, 2014

Contents

1	Introduction - Open Research Methods	1
2	What is R?	2
2.1	R user-interfaces	3
2.2	Structure of R-project	3
3	Learning the language	4
4	Use of R language	5
5	Conclusions	5
6	References	7

1 Introduction - Open Research Methods

The debate on open science in the context of Social sciences and humanities (SSH) has been predominantly focusing on open access to research publication and opening up the various types of digital research data (open research data). The openness of research methods has received a lot less attention.

I can think of at least two reasons for that. On the one hand, research methods in SSH have been predominantly qualitative where software has played only a supporting role. Such research methods, let's take *discourse analysis*, have always been open. On the other hand, the quantitative fields of SSH have mostly

used survey, register, statistics and other, often closed, tailor-made data that custom statistical tools such as SPSS, Stata or Excel are well suited for. However, the future of SHH looks somewhat different as the quantity and multiplicity of sources of digital data are challenging both traditional approaches in SSH, the purely qualitative approach and custom tools approach in quantitative analysis. As Gary King (2014) describes it:

An important driver of the change sweeping the field is the enormous quantities of highly informative data inundating almost every area we study. In the last half-century, the information base of social science research has primarily come from three sources: survey research, end-of-period government statistics, and one-off studies of particular people, places, or events. In the next half-century, these sources will still be used and improved, but the number and diversity of other sources of information are increasing exponentially and are already many orders of magnitude more informative than ever before.

As the role of software and computation in the research process is becoming more important, the questions of licensing, ownership, modification and distribution of that software become increasingly important, too. This chapter will introduce one viable option for analysing your data with called R.

2 What is R?

R is one of the most popular platforms for data analysis and visualization currently available. R is distributed under the terms of the *GNU General Public License* so it is free and open source and it can be distributed under those conditions. R is available from [Comprehensive R Archive Network](#) (CRAN). The name R comes from the first names of two New Zealand statisticians Ross Ihaka and Robert Gentleman who created the language in late 1990's.

R can be regarded as an implementation of the S language which was developed at Bell Laboratories in 70's by Rick Becker, John Chambers and Allan Wilks (Venables, Smith, and Team 2013). R is an object-oriented programming language which means that unlike in SPSS or SAS that give you abundant information on particular model you implement, R only creates a *fit object* in memory that can be used in subsequent analysis. This structure of R directs the user to implement the data-analysis as stepwise process which becomes very useful later on when solving complex research problems using vast and messy data typical for emerging computational SSH research.

2.1 R user-interfaces

R runs in Windows, Mac OS X and GNU/Linux operating systems on local computer, but a different server implementations are becoming increasingly popular, such as [R-Fiddle](#) or [rnotebook](#). The most basic user interface for R is *console*, which allows user to type in commands and outputs the results of the analysis. If the results is a plot a pop-up graphical windows is opened. There are several graphical user interfaces (GUI) in R that may be helpful in the beginning, like [RCommander](#) or [Deducer](#). Perhaps the most productive way for using R is through a [integrated development environment](#) (IDE) that provide the user, in addition to console, several useful functionalities for controlling the whole research project. [RStudio](#) has gained a lot of popularity in last couple of years and is also my personal favourite IDE. It combines the console with script editor, plot browser, file browser and environment window. If the user uses plain text (latex or markdown) for typesetting the texts, RStudio has a tailored text editor and support for version control either in [git](#) or in [subversion](#). In addition, RStudio has native support for html-based presentation graphics using [reveal.js](#)-framework. All these operations makes it possible to squeeze the whole research process within a single software environments from planning to publishing. Rstudio can also be run on a remote server through web browser. The RStudio company has another exciting open source tool for R called [shiny](#) that can be used for creating interactive web applications such as this [experimental gadget of mine](#).

2.2 Structure of R-project

For someone new to R, the peculiar structure of the language creates a very steep learning curve. Same applies to learning how the whole project is organised.

Official name [The R Project for Statistical Computing](#) refers both to the centrally maintained core as well as R's distributed structure of contributed extensions, called *packages*. Packages in R are collections of functions and/or data that are packaged for convinieny. Installing a package broadens your the functionality of your R installation. Basic R installation consists of so called *base installation* that includes the core with some 25 additional *packages* for the most basic functionality. The core of the language is maintained by *R Development Core Team*, but the additional packages are developed and maintained by individual developers and research institutes. R users often create packages for themselves, but if one thinks the package could be useful for other users too, the packages can be distributed through repositories.

CRAN is the “official” repository for contributed packages and currently hosts 5150 packages that can be used to extend R. In the last couple of years various code hosting sites such as [GitHub](#) have become increasingly important resources especially for collaborative development of new packages. Github hosts currently roughly 1500 packages for R. [Bioconductor](#) is another separate package repository,

but can be regarded as *domain specific* for it hosts packages *for the analysis and comprehension of high-throughput genomic data*. Other such domain specific projects are for example [rOpenSci](#) and emerging [rOpenGov](#) that provide tools for open science and open government data, respectively.

3 Learning the language

As the internet has brought together the vast community around R and internet has become the main channel for delivering instructions for R. The official *Introduction to R* by Venables, Smith, and Team (2013) is an important document to master when getting into the language. Besides this general introduction R-project has also a domain specific structure where you can start learning from so called [task views](#). For SSH researcher the [social sciences](#) and [natural language processing](#) task views are good places to begin with.

Discussions and announcement on R happen mainly through [R official mailing lists](#) that have their own lists for development and user help. [R help](#) is the main list for general help and receives tens of mails per day. Most of the individual packages have their own mailing list for development where anyone can join if wanting to contribute to the packages.

The official mailing lists have recently been challenged by so called Question & Answer -sites like [Stack Overflow](#) in delivering solutions for one-off user questions. Stack Overflow has currently almost 47000 questions tagged with R. In comparison to proprietary software, there are 2014 questions tagged with SAS, 616 with Stata and 362 SPSS. These figures are used as one indicator of the [increasing popularity of R](#). Besides the Question & Answer sites, there are hundreds of blogs discussing specific analytical problem using R and feeds from the blogs are aggregated in [R-bloggers](#)-website.

Another, more formal channel for distributing and communicating R have become the so called massive open online courses (MOOC). MOOCs seem to work well for teaching programming and many courses in [Coursera](#) and [EdX](#) have become hugely popular attracting tens of thousands of students each year. The free licensing of R has made it primary language on these courses as it is basically the only viable alternative for teaching statistical programming for massive crowds.

Aside with vibrant internet community more and more books are being published on R. Books can be put in three categories. First are the general introductions to statistics using R. *Discovering Statistics Using R* by A. Field, Miles, and Field (2012) and *R in Action: Data Analysis and Graphics With R* by Robert Kabacoff (2013) are popular examples of that category. Second there are more and more books addressing how to solve some specific analytical problems using R. A prime examples of books in this category are *Complex Surveys: A Guide to Analysis Using R* by Thomas Lumley (2011), *Text Analysis with R for Students of Literature* by Matthew Jockers (forthcoming), *R Graphics Cookbook* by Chang

(2012) and *Dynamic documents with R and knitr* by Yihui Xie (2014). A third category are the books that focus on specific theoretical issue in statistics and use R as a primary language to demonstrate this. Such books are for instance *Bayesian Data Analysis* by Gelman et al. (2013) or *Multilevel Analysis: An Introduction to Basic and Advanced Multilevel Modeling* by Snijders and Bosker (2011).

4 Use of R language

Throughout it's existence the main use of R has been implementation of new statistical methods. This is still the case and implementations of new statistical methods are usually first available in R. However, applied statistics have become more and more important area of expansion as researchers from fields that apply statistical methods have started to migrate into R. *Bioconductor* was already mentioned as an example of a domain specific initiative to apply R in their analysis, for genome data in this case. Natural sciences in general have been early adopters and for example in Geographical Information Systems (GIS) the R has started rival proprietary GIS-software. In this case for example it is possible to combine traditional statistical methods and geographic data in one environment. In SSH this may be of useful as there are a lot of spatial data available and you may want run cluster thematically, but also visualize as maps. Matthew Jockers (forthcoming) book is one of the few attempts to foster use of R in humanities. In digital humanities blogosphere there are few others besides [Jockers blog](#) that are worth reading, namely [W. Caleb McDaniel from Rice University](#) and [Quantifying Memory blog](#) by Rolf Fredheim from Cambridge.

Addition to academic applications R has become a major player in business analytics. This is largely due to companies as [Revolution Analytics](#) that have started providing consultation and creating proprietary application for enterprise needs. Annual [R/Finance 2014](#)-conference gives nice overview of adoption of R in banking and insurance sectors. For example [Google](#) uses R in-house and also provides packages as [r-google-analytics](#) or [rgooglevis](#). One emerging field is so called data journalism where major players as [New York Times](#) or [Guardian](#)

5 Conclusions

R is not the only alternative for proprietary software. [Python](#) is another viable option especially for someone looking for more general purpose language that also masters data analysis. [Whether python is going to displace R](#) has recently been debated in data science blogosphere as the data is becoming mainstream, but R is still remaining hard to learn and programmers rather want to extend the language they already know rather than learn a new one. For very intensive computation [Julia](#) is becoming popular open source option, too. It is still in

early phase of development, but is already viable option if processing time is important.

However, for scientific work I would emphasize the licensing of the software more than the name of the technology. It may be that the recent buzz around digital data in SSH marks only the beginning of the data intensive research and for someone wanting to success in that game it will be equally important to develop the substantial understanding as well as technological understanding to get grip of the new emerging tools. R is a prime example of this development where academics have taken a major role in software development and created tools that suit better for many contemporary research problems than proprietary software. This development may well go on and therefore it is advisable for someone who is interested in learning this techniques to carefully look at the licensing before investing time and effort in learning the tools. Free and open source tools are great in this respect as once you pick up the skills to use the product, you will soon find that it should be improved for your purposes. You can learn how the code works, write improvements and then publish them for the community. In addition, free licensing also allows you to teach the use of the software and distribute it. Open source research software is not always the easiest and quickest way to get job done, but in the long run they are often worth the time invested.

Finally, the openness of the computation research methods is important from the reproducibility of your research. Along with demands for open access of research publications there are tendencies that more and more journals in computational sciences will require both the data and algorithms behind the results to published together with the article. As SSH scholars are moving towards computational analysis this issue of reproducibility should also be considered in choosing tools. R is a great tool that fullfills all these conditions, but are several other out there, too. After all, it is not necessary for all to become software developers, but to have basic understanding and pair with developers who know more.

Steve Lohr (2013) interviewed some leading digital humanists in New York Times article *Literary History, Seen Through Big Data's Lens* on the future of social sciences posed a question whether these emerging computational technologies will undermine the role qualitative research in SHH. Matthew Jocker emphasized that finding the right questions and flaws in the analysis still require deep understanding of the field:

“You’ll always need both,” says Mr. Jockers, the literary quant. “But we’re at a moment now when there is much greater acceptance of these methods than in the past. There will come a time when this kind of analysis is just part of the tool kit in the humanities, as in every other discipline.”

Quantitative tools in the humanities and the social sciences, as in other fields, are most powerful when they are controlled by an

intelligent human. Experts with deep knowledge of a subject are needed to ask the right questions and to recognize the shortcomings of statistical models.

Also professor Gary (King 2014) from Harvard claims that this development is requires new kind of cooperation between scholars and fields of research as the analysis of large digital data requires skills that can't be found from traditional fields of social sciences.

Through collaboration across fields, however, we can begin to address the interdisciplinary substantive knowledge needed, along with the engineering, computational, ethical, and informatics challenges before us. (King 2014)

In addition, King (2014) assumes that this deveploment will blur the dichotomy between qualitative and quantitative analysis and he sees a future where both traditions have merged into social sciences where the old problems are solved in collaboration.

Instead of quantitative researchers trying to build fully automated methods and qualitative researchers trying to make due with traditional human-only methods, both now are heading toward, using, or developing computer-assisted methods that empower both groups. This development has the potential to end the divide, to get us working together to solve common problems, and to greatly strengthen the research output of social science as a whole.

6 References

- Chang, Winston. 2012. *R Graphics Cookbook*. O'Reilly.
- Field, Andy, Jeremy Miles, and Zoë Field. 2012. *Discovering Statistics Using R*. SAGE.
- Gelman, Andrew, John B. Carlin, Hal S. Stern, David B. Dunson, Aki Vehtari, and Donald B. Rubin. 2013. *Bayesian Data Analysis, Third Edition*. CRC Press.
- Jockers, Matthew. forthcoming. *Text Analysis with R for Students of Literature*. Quantitative Methods in the Humanities and Social Sciences. Springer.
- Kabacoff, Robert. 2013. *R in Action: Data Analysis and Graphics With R*. MANNING PUBN.
- King, Gary. 2014. "Restructuring the Social Sciences: Reflections from Harvards Institute for Quantitative Social Science." *PS: Political Science and Politics* 47: 165–172. http://journals.cambridge.org/repo_A9100Nlq.

- Lohr, Steve. 2013. “Literary History, Seen Through Big Data’s Lens.” *The New York Times* (January). <http://www.nytimes.com/2013/01/27/technology/literary-history-seen-through-big-datas-lens.html>.
- Lumley, Thomas. 2011. *Complex Surveys: A Guide to Analysis Using R*. John Wiley & Sons.
- Snijders, Tom A. B., and Roel Bosker. 2011. *Multilevel Analysis: An Introduction to Basic and Advanced Multilevel Modeling*. Second Edition. Sage Publications Ltd.
- Venables, William N., David M. Smith, and R. Development Core Team. 2013. *An Introduction to R*. Network Theory. <http://www.math.vu.nl/sto/onderwijs/statlearn/R-Binder.pdf>.
- Xie, Yihui. 2014. *Dynamic Documents with R and Knitr*. CRC Press.