

Use R!

Douglas A. Luke

A User's Guide to Network Analysis in R



Springer

Use R!

Series Editors:

Robert Gentleman Kurt Hornik Giovanni Parmigiani

More information about this series at <http://www.springer.com/series/6991>

Use R!

Albert: Bayesian Computation with R (2nd ed. 2009)

Bivand/Pebesma/Gómez-Rubio: Applied Spatial Data Analysis with R (2nd ed. 2013)

Cook/Swayne: Interactive and Dynamic Graphics for Data Analysis:
With R and GGobi

Hahne/Huber/Gentleman/Falcon: Bioconductor Case Studies

Paradis: Analysis of Phylogenetics and Evolution with R (2nd ed. 2012)

Pfaff: Analysis of Integrated and Cointegrated Time Series with R (2nd ed. 2008)

Sarkar: Lattice: Multivariate Data Visualization with R

Spector: Data Manipulation with R

Douglas A. Luke

A User's Guide to Network Analysis in R



Springer

Douglas A. Luke
Center for Public Health Systems Science
George Warren Brown School
of Social Work
Washington University
St. Louis, MO, USA

ISSN 2197-5736

Use R!

ISBN 978-3-319-23882-1

DOI 10.1007/978-3-319-23883-8

ISSN 2197-5744 (electronic)

ISBN 978-3-319-23883-8 (eBook)

Library of Congress Control Number: 2015955739

Springer Cham Heidelberg New York Dordrecht London

© Springer International Publishing Switzerland 2015

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made.

Printed on acid-free paper

Springer International Publishing AG Switzerland is part of Springer Science+Business Media (www.springer.com)

*To my most important social network—Sue,
Alina, and Andrew*

Preface

In early 2000, Stephen Hawking said that “...the next century will be the century of complexity.” If his prediction is true, the implication is that we will need new scientific theories, data collection methods, and analytic techniques that are appropriate for the study of complex systems and behavior. Network science is one such approach that views the world through a network lens, where physical and social systems are made up of heterogeneous actors who are connected to one another through different types of relational ties. Network analysis is the set of analytic tools used to study these types of systems. Over the past several decades network analysis has become an increasingly important part of the analytic toolbox for social, health, and physical scientists.

Until recently, network analysis required specialized software, both for network data management and analyses. However, starting around 2000, network analytic tools became available in the R statistical programming environment. This not only made network analytic techniques more visible to the broader statistical community but also provided the breadth and power of R’s data management, graphic visualization, and general statistical modeling capabilities to the network analyst community.

As the title suggests, this book is a user’s guide to network analysis in R. It provides a practical hands-on tour of the major network analytic tasks that can currently be done in R. The book concentrates on four primary tasks that a network analyst typically concerns herself with: network data management, network visualization, network description, and network modeling. The book includes all the R code that is used in the network analysis examples. It also comes with a set of network datasets that are used throughout the book. (See Chap. 1 for more details on the structure of the book, as well as instructions on how to obtain the network data.) The book is written for anybody who has an interest in doing network analysis in R. It can be used as a secondary text in a network science or analysis class or can simply serve as a reference for network techniques in R.

This book would not exist without the help, support, guidance, and mentoring I have received over the last 30 years from my own personal and professional social networks. In the mid-1980s I took a graduate network analysis class from Stan Wasserman at the University of Illinois in Champaign. I remember being excited

about this new way to analyze data, but thought that I was not likely to ever use it in my career. However, my colleagues in psychology and public health encouraged me in my early work exploring how network analysis could answer important research and evaluation questions. These include Julian Rappaport, Ed Seidman, Bruce Rapkin, Kurt Ribisl, Sharon Homan, Ross Brownson, and Matt Kreuter. Whether they know it or not, I have been inspired and encouraged by an amazing group of network and systems scientists, including Tom Valente, Steve Borgatti, Martina Morris, Tom Snijders, Scott Leischow, Patty Mabry, Stephen Marcus, and Ross Hammond. My best network ideas have come from my friends and colleagues at the Center for Public Health Systems Science, particularly Bobbi Carothers, Amar Dhand, Chris Robichaux, and Nancy Mueller. I am especially grateful to the students in my network analysis classes and workshops over the years; they have not only improved this book, but they have improved my thinking about network analysis. A very special thank you to Jenine Harris. Jenine was my first doctoral student, now I am inspired by the rigor and elegance of her own work in network science. I would also like to thank the Centers for Disease Control and Prevention, the National Institutes of Health, and the Missouri Foundation for Health for providing research and evaluation support that allowed me to develop and refine my approach to network analysis. Finally, my deepest thanks go to my family. They gave me specific suggestions about the content, provided me space and time to work hard on this book (including a crucial Father's Day gift), and cheered me on when I most needed it. Thank you, Sue, Ali, and Andrew.

St. Louis, MO, USA
July, 2015

Douglas A. Luke

Contents

1	Introducing Network Analysis in R	1
1.1	What Are Networks?	1
1.2	What Is Network Analysis?	3
1.3	Five Good Reasons to Do Network Analysis in R	4
1.3.1	Scope of R	4
1.3.2	Free and Open Nature of R	5
1.3.3	Data and Project Management Capabilities of R	5
1.3.4	Breadth of Network Packages in R	6
1.3.5	Strength of Network Modeling in R	6
1.4	Scope of Book and Resources	6
1.4.1	Scope	6
1.4.2	Book Roadmap	7
1.4.3	Resources	8
Part I Network Analysis Fundamentals		
2	The Network Analysis ‘Five-Number Summary’	11
2.1	Network Analysis in R: Where to Start	11
2.2	Preparation	11
2.3	Simple Visualization	12
2.4	Basic Description	12
2.4.1	Size	12
2.4.2	Density	14
2.4.3	Components	15
2.4.4	Diameter	15
2.5	Clustering Coefficient	16
3	Network Data Management in R	17
3.1	Network Data Concepts	17
3.1.1	Network Data Structures	17
3.1.2	Information Stored in Network Objects	20

3.2	Creating and Managing Network Objects in R	21
3.2.1	Creating a Network Object in <code>statnet</code>	21
3.2.2	Managing Node and Tie Attributes	24
3.2.3	Creating a Network Object in <code>igraph</code>	28
3.2.4	Going Back and Forth Between <code>statnet</code> and <code>igraph</code> ..	30
3.3	Importing Network Data	30
3.4	Common Network Data Tasks	32
3.4.1	Filtering Networks Based on Vertex or Edge Attribute Values	32
3.4.2	Transforming a Directed Network to a Non-directed Network	39

Part II Visualization

4	Basic Network Plotting and Layout	45
4.1	The Challenge of Network Visualization	45
4.2	The Aesthetics of Network Layouts	47
4.3	Basic Plotting Algorithms and Methods	49
4.3.1	Finer Control Over Network Layout	50
4.3.2	Network Graph Layouts Using <code>igraph</code>	52
5	Effective Network Graphic Design	55
5.1	Basic Principles	55
5.2	Design Elements	55
5.2.1	Node Color	56
5.2.2	Node Shape	60
5.2.3	Node Size	62
5.2.4	Node Label	66
5.2.5	Edge Width	68
5.2.6	Edge Color	69
5.2.7	Edge Type	70
5.2.8	Legends	71
6	Advanced Network Graphics	73
6.1	Interactive Network Graphics	73
6.1.1	Simple Interactive Networks in <code>igraph</code>	74
6.1.2	Publishing Web-Based Interactive Network Diagrams	74
6.1.3	Statnet Web: Interactive <code>statnet</code> with <code>shiny</code>	77
6.2	Specialized Network Diagrams	77
6.2.1	Arc Diagrams	78
6.2.2	Chord Diagrams	79
6.2.3	Heatmaps for Network Data	82
6.3	Creating Network Diagrams with Other R Packages	84
6.3.1	Network Diagrams with <code>ggplot2</code>	84

Part III Description and Analysis

7 Actor Prominence	91
7.1 Introduction	91
7.2 Centrality: Prominence for Undirected Networks	92
7.2.1 Three Common Measures of Centrality	93
7.2.2 Centrality Measures in R	95
7.2.3 Centralization: Network Level Indices of Centrality	96
7.2.4 Reporting Centrality	97
7.3 Cutpoints and Bridges	101
8 Subgroups	105
8.1 Introduction	105
8.2 Social Cohesion	106
8.2.1 Cliques	107
8.2.2 k-Cores	110
8.3 Community Detection	115
8.3.1 Modularity	115
8.3.2 Community Detection Algorithms	118
9 Affiliation Networks	125
9.1 Defining Affiliation Networks	125
9.1.1 Affiliations as 2-Mode Networks	126
9.1.2 Bipartite Graphs	126
9.2 Affiliation Network Basics	127
9.2.1 Creating Affiliation Networks from Incidence Matrices	127
9.2.2 Creating Affiliation Networks from Edge Lists	129
9.2.3 Plotting Affiliation Networks	130
9.2.4 Projections	131
9.3 Example: Hollywood Actors as an Affiliation Network	133
9.3.1 Analysis of Entire Hollywood Affiliation Network	134
9.3.2 Analysis of the Actor and Movie Projections	139

Part IV Modeling

10 Random Network Models	147
10.1 The Role of Network Models	147
10.2 Models of Network Structure and Formation	148
10.2.1 Erdős-Rényi Random Graph Model	148
10.2.2 Small-World Model	151
10.2.3 Scale-Free Models	154
10.3 Comparing Random Models to Empirical Networks	160

11 Statistical Network Models	163
11.1 Introduction	163
11.2 Building Exponential Random Graph Models	165
11.2.1 Building a Null Model	167
11.2.2 Including Node Attributes	169
11.2.3 Including Dyadic Predictors	171
11.2.4 Including Relational Terms (Network Predictors)	175
11.2.5 Including Local Structural Predictors (Dyad Dependency) ..	177
11.3 Examining Exponential Random Graph Models	179
11.3.1 Model Interpretation	179
11.3.2 Model Fit	180
11.3.3 Model Diagnostics	183
11.3.4 Simulating Networks Based on Fit Model	183
12 Dynamic Network Models	189
12.1 Introduction	189
12.1.1 Dynamic Networks	189
12.1.2 RSiena	191
12.2 Data Preparation	192
12.3 Model Specification and Estimation	198
12.3.1 Specification of Model Effects	198
12.3.2 Model Estimation	203
12.4 Model Exploration	203
12.4.1 Model Interpretation	203
12.4.2 Goodness-of-Fit	209
12.4.3 Model Simulations	212
13 Simulations	217
13.1 Simulations of Network Dynamics	217
13.1.1 Simulating Social Selection	218
13.1.2 Simulating Social Influence	228
References	235

Chapter 1

Introducing Network Analysis in R

Begin at the beginning, the King said, very gravely, and go on till you come to the end: then stop. (Lewis Carroll, Alice in Wonderland)

1.1 What Are Networks?

This book is a user’s guide for conducting network analysis in the R statistical programming language. Networks are all around us. Humans naturally organize themselves in networked systems. Our families and friends form personal social networks around each of us. Neighborhoods and communities organize themselves in networked coalitions to advocate for change. Businesses work with (and against) each other in complex, interlocking networks of trade and financial partnerships. Public health is advanced through partnerships and coalitions of governmental and NGO organizations (Luke and Harris 2007). Nations are connected to one another through systems of migration, trade, and treaty obligations.

Moreover, non-human networks exist almost anywhere you look. Our genes and proteins interact with one another through complex biological networks. The human brain is now viewed as a complex network, or ‘connectome’ (Sporns 2012). Similarly, human diseases and their underlying genetic roots are connected as a ‘diseasome’ (Barabási 2007). Animal species interact in many complex ways, one of which is a networked food-web that describes interactions in ‘who-eats-whom’ relationships. Information itself is networked. Our legal system is built on an inter-connecting network of prior legal decisions and precedents. Social and scientific progress is driven by a diffusion of innovation process by which information is disseminated across connected social systems, whether they are Iowa corn farmers (Rogers 2003) or public health scientists (Harris and Luke 2009). It appears that one of the ways the universe is organized is with networks.

So what is a network? Figures 1.1 and 1.2 present two examples of important and interesting social networks. Figure 1.1 presents the contact network of the 19 9–11 hijackers, based on the work of Valdis Krebs (2002). Every social network is made up of a set of actors (also called nodes) that are connected to one another via some type of social relationship (also called a tie). In the figure, nodes are the circles and the ties are the lines connecting some of the nodes. The network shows

us that the hijackers had some contact with one another before September 11th, but the network is not very densely connected and there appears to be no prominent network member who is connected to all or even most of the other hijackers.

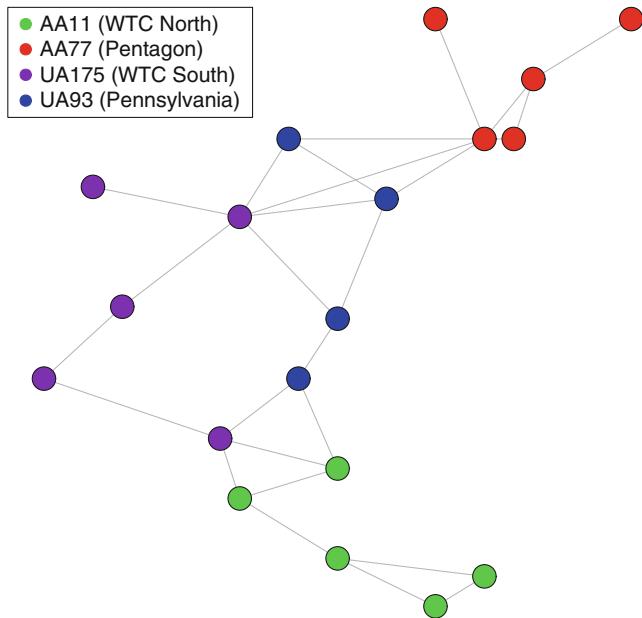


Fig. 1.1 Network of 9–11 hijackers

The second example in Fig. 1.2 is from a very different sort of social network. Here the nodes are members of the 2010 Netherlands FIFA World Cup team, who went on to lose in the final to Spain. The ties represent passes between the different players during the World Cup matches. The arrows show the directional pattern of the passes. We can see that the goalkeeper passed primarily to the defenders, and the forwards received passes primarily from the midfielders (except for #6, who appears to have a different passing pattern than the other two forwards).

These two examples may appear to have little in common. However, they both share a fundamental characteristic common to all social networks. The social patterns that are displayed in the network figures are not random. They reflect underlying social processes that can be explored using network science theories and methods. The terrorist network has no prominent leader and is not tightly interconnected because it makes the network harder to detect or disrupt. The pattern of passing ties in the soccer network reflects the assigned positions of the players, the rules of the game, and the strategies of the coach. The network analysis does not ‘know’ about any of those rules or strategies. Yet, network analysis can be used to reveal these patterns that reflect the underlying rules and regularities.