

Lecture Notes in Social Networks

Katharina A. Zweig

Network Analysis Literacy

A Practical Approach to the Analysis of
Networks

 Springer

Lecture Notes in Social Networks

Series editors

Reda Alhadj, University of Calgary, Calgary, AB, Canada

Uwe Glässer, Simon Fraser University, Burnaby, BC, Canada

Advisory Board

Charu Aggarwal, IBM T.J. Watson Research Center, Hawthorne, NY, USA

Patricia L. Brantingham, Simon Fraser University, Burnaby, BC, Canada

Thilo Gross, University of Bristol, Bristol, UK

Jiawei Han, University of Illinois at Urbana-Champaign, IL, USA

Huan Liu, Arizona State University, Tempe, AZ, USA

Raúl Manásevich, University of Chile, Santiago, Chile

Anthony J. Masys, Centre for Security Science, Ottawa, ON, Canada

Carlo Morselli, University of Montreal, QC, Canada

Rafael Wittek, University of Groningen, The Netherlands

Daniel Zeng, The University of Arizona, Tucson, AZ, USA

More information about this series at <http://www.springer.com/series/8768>

Katharina A. Zweig

Network Analysis Literacy

A Practical Approach to the Analysis
of Networks

 Springer

Katharina A. Zweig
TU Kaiserslautern
FB Computer Science
Graph Theory and Analysis of Complex
Networks
Kaiserslautern
Germany

ISSN 2190-5428 ISSN 2190-5436 (electronic)
Lecture Notes in Social Networks
ISBN 978-3-7091-0740-9 ISBN 978-3-7091-0741-6 (eBook)
DOI 10.1007/978-3-7091-0741-6

Library of Congress Control Number: 2016948283

© Springer-Verlag GmbH Austria 2016

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

This Springer imprint is published by Springer Nature
The registered company is Springer-Verlag GmbH Austria
The registered company address is: Prinz-Eugen-Strasse 8-10, 1040 Wien, Austria

*To Winfried Zweig
Thanks for all the good discussions
on network analysis.*

Foreword

This is a delightful book. It 'is so easy to read, and you can almost accidentally learn quite a bit of network science without even noticing it. Written in a playful manner, it tends to enliven the brain rather than put it to sleep—quite a change from the usual pedantic tone. It's a quirky book that does not try to be systematic. For example, it does not cover “community detection” (that's cluster analysis to you social scientists). As a result, the book has a great deal of personality.

But what I really like about the book is the subtext. What it's actually about, in my opinion, is how to think, and here, that means how to think with models. Most academics are very gullible when it comes to concepts outside their disciplines. Within their area, any new idea or phrasing is treated with withering skepticism, but outside their area, they adopt ideas with the speed of teenagers adopting slang or fashion. Thus, a management scholar hears about small worlds and clustering coefficients and immediately shoehorns them into their next study. A physicist learns about betweenness centrality, and suddenly, there are 500 papers that reference the idea. If the first paper associates betweenness with influential spreaders in the spread of a disease, all of the following papers do the same. If you internalize this book, you won't make that mistake. You will realize that although there is a sense in which network measures are tools like hammers, there is much more to them. Hammers work pretty much the way they work in any setting, but using a network measure implicitly entails fitting a model of how things work. And if the model doesn't fit, the measure doesn't either.

Curiously, although I associate model-based thinking with the physical sciences, my experience is that both physical and social scientists are equally likely to have this mindless, “pluginski” attitude about network concepts. Therefore, I think this book would be useful for both audiences. But since the content of the book is mostly drawn from what Katharina calls the “network science” field (as opposed to

the “social network analysis” field), I’m guessing it will appeal mostly to budding physical scientists. Too bad, because if there was ever an introduction to network science that was especially suitable for social scientists, this is it.

I look forward to seeing this in print.

Steve Borgatti
Lexington KY, USA

Preface

On the other hand, this rapid expansion [of complex network science] creates the risk that existing methods may be misapplied or misinterpreted, leading to inappropriate conclusions and generally poor results. (Carter Butts: “Revisiting the Foundations of Network Analysis,” Science, 325, 414–416, 2009)

After finishing a study of biochemistry and in the middle of a bioinformatics study, I started my work as a doctoral student in the field of “Algorithm design and computational complexity” in 2003. I was immediately attracted to a budding new field, *complex network science*, which had just taken off some years ago.

I was lucky enough to meet Ulrik Brandes early in this endeavor, and he invited me to participate in the now classic textbook edited by him and Thomas Erlebach: “Network analysis—Methodological Foundations.” With other doctoral students, I was assigned to the chapters on *centrality indices*. In the beginning, I was overwhelmed by the dozens of different indices that had been proposed so far and the seemingly never-ending flow of newly proposed centrality indices. The argumentation almost always went along the following lines: “So far, these indices have been proposed. In the new data set X , none of these measures matches with the intuition. Thus, we propose the new measure Y that matches our intuition of centrality in this network.” I was lost which index to take in any specific situation.

Finding a first online version of Stephen P. Borgatti’s paper on “Centrality and network flow” was a revelation: Borgatti basically says that a centrality index is a predictor of which node is used most heavily in a given network flow or network process. While others like Freeman had also hinted at a relation between processes on a network and a measure to quantify the network’s structure, Borgatti was the first to make a tight connection between a process of interest and the measure to quantify the indirect effects induced by this network process. He also stated quite clearly that a mismatch between a complex network, the network process of interest, and the centrality index will lead to uninterpretable results: “the off-the-shelf formulas for centrality measures are fully applicable only for the

specific flow processes they are designed for, and (...) when they are applied to other flow processes they get the ‘wrong’ answer.”¹

This book is based on the idea that network processes and network analytic measures are even more intertwined, beyond the set of centrality indices. In the last ten years I have generalized this idea to all kinds of distance—and walk-based measures. The main hypothesis of this book is as follows:

Note 1. To interpret the values of a distance-based measure, the way of calculating the distance must be matched to the process of interest. To interpret any walk-based measure, the set of walks used by the measure needs to be closely adapted to the process.

This includes the whole process of data observation, preprocessing, representation as a network, stating a network process of interest, and choosing a network analytic method to analyze it. It is the book that I would have loved to have at the beginning of my doctoral research.

Intended Audience

There seem to be three types of groups pursuing network analytic projects:

1. Groups consisting of scientists with a heap of data that want to analyze their data by network analytic methods—henceforth called *data experts*.
2. Groups consisting of scientists that primarily devise network analytic methods and then search for data that can be analyzed by their newly devised method—henceforth called *method experts*.
3. A quite small set of interdisciplinary groups, consisting of data **and** method experts.

As a biochemist, I was clearly in the first group as a *data expert* that was overwhelmed by the choice of methods; later, as an algorithm designer and method expert, it became clear that applying the best and most beautiful method to data and a research question it does not match with, is not helpful either.

This book will stress that people from both groups need to be *literate* in the other group’s regime: If a data expert creates a beautiful data set which can be represented and analyzed as a network, it is important not to miss any vital pattern just because a particularly suitable method is not known to him or her. Similarly, for any method expert, it is vital to understand the data to which a chosen method is applied. In particular, it is not enough to just reference to the data expert’s publication and to roughly know what the vertices and edges represent, but it is necessary to understand in detail how the data were produced, to know the odds of observing false-positive and false-negative relationships, and to know whether the resulting network is complete or not. However, many data experts do not include this

¹Stephen P. Borgatti: “Centrality and Network Flow”, *Social Networks* 27, 55–71, 2005.

information in their publications, for example, because the community from which the data originate is well aware of the applied procedures.

This book tries to build the bridge between the two groups and to show the different perspectives they have on their subjects and projects.

The Ideal Reader

Yes, I have some requirements toward you, my dear reader. Obviously, you are a data expert who thinks that network analysis could be helpful to reveal the most exciting mysteries in your field—and so do I. With this book, I will equip you with the necessary questions you need to ask your method expert to understand whether your research question matches with his or her method.

Or you are a method expert, maybe a mathematician or a computer scientist, and your advisor just gave you this piece of data and asked you to design a method to analyze it—then, this is the book that will help you to understand which questions to ask your data provider.

It is just the book I wanted to have when I was about one year into my doctoral studies, still overwhelmed by the amazing flexibility of network analysis and underwhelmed by the number of good guidelines to use it: guidelines on how to actually represent a complex system by a network or how to choose the best method to analyze it, and how all of this is connected to my research question. I was baffled by the daring approach of physicists that simplified complex systems beyond recognition to a set of nodes and edges—and at the same time, I was intrigued by the potential of this new approach. However, I was still biochemist enough to wonder whether there is actually a line where simplification needs to stop, in order to find contextually meaningful results. If you are at this point in your career, I wrote this book for you.

If you are not quite there yet, you might want to read the very good collection of papers, edited by Mark Newman, Albert-László Barabási, and Duncan J. Watts, called “The structure and dynamics of networks.”² For those in a hurry, the following papers are the minimally required prerequisites to get a feeling for the field:

1. Start with the two papers that opened the field of (social) network analysis to a much broader community and transformed it into complex network analysis: The first was published 1998 by Duncan Watts and Steven H. Strogatz under the title “Collective dynamics of small-world networks,” *Nature* 393, pp. 440–442. The paper introduced the small-world model. The second influential paper was published in 1999 by Albert-László Barabási and Réka Albert under the title “Emergence of scaling in random networks” in *Science* 286, pp. 509–512. It introduced the notion of scale-free networks and a model to produce them, the preferential attachment model. Both papers are shortly summarized in Chap. 6.

²Published by Princeton University Press, Princeton and Oxford, 2006

2. The first paper by Barabási was quickly followed by a disturbing one which showed that scale-free networks built with the preferential attachment model are robust against random failures of nodes but very sensitive to attacks on their most connected nodes: Albert, Jeong, and Barabási published these findings in 2000 under the title “Error and attack tolerance of complex networks” in *Nature* 406, pp. 378–382.
3. For the course of this book, the work on so-called *network motifs* by Uri Alon’s group is very much important.³ Other disciplines had already started earlier to explicitly compare a structural value found in a network with the expected one, for example ecology.⁴ For the field of complex network analysis, the articles by Alon et al. were the first widely visible ones that proposed to assign a significance value to observed results by comparing the observation with the expectation.
4. The articles above are written by physicists. Read now the view of the sociologists, as stated by Borgatti et al.’s paper on “Network Analysis in the Social Sciences,” published in *Science*, 323, pp. 892–895, in 2009.

By reading these papers, you might notice that the publications in the field of complex network analysis come from very different publication venues. This is caused by the very interdisciplinary origin of the field. For example, for computer scientists, it is common to publish their original research in a conference proceeding, and some of their conferences are as reliable and respected as journals. For physicists, a conference is a place to meet and exchange ideas, but most often, they report recent work that was already published elsewhere. Please read the chapter in the appendix discussing different publication styles, where to find which information and how to differentiate peer-reviewed from unreviewed publications.

Now, you are well-prepared for an instruction on how to read this book!

How to Read This Book

As a reader, I mostly skip these sections on “How to read this book,” so I make it extra-short: This is a book to be read from left to right and from top to bottom, or to dip in as you please. The exercises are intended to deepen the understanding of the methods introduced in the text. Moreover, they teach what questions to ask whenever you make acquaintance with a new measure. Almost all exercises can be solved on two levels: a verbal, explanatory solution with the help of an example and

³I suggest to read: Ron Milo et al.: “Network Motifs: Simple Building Blocks of Complex Networks,” *Science* 298, pp. 824–827, 2002 and Ron Milo et al.: “Superfamilies of Evolved and Designed Networks,” *Science* 303, pp. 1538–1542, 2004.

⁴Nicholas J. Gotelli and Gary R. Graves: “Null-Models in Ecology,” Smithsonian Institution Press, 1996.

by proof. For courses with mathematicians, physicists, and computer scientists, I normally require a proof for these exercises.

The book is divided in three parts: Part I gives you an overview of the field and names the necessary definitions. Part II is devoted to the most important methods, starting with some classic, network analytic measures, a basic discussion on how to represent data as complex networks, various random graph models and their use in network analysis, and centrality indices. Most importantly, it also contains a chapter on how to analyze a measure that you encounter somewhere. Both parts are just the preparation for the core of this book, Part III, which describes various aspects of *network analysis literacy*: when data cannot be represented as a network, when a method's results are difficult to interpret, and finally, why network analysis is a field that even sometimes requires an ethical perspective.

So, where would I recommend you to start? When you are an absolute beginner in network analysis, start with—surprise—Chap. 1. If you are already confused by network analysis, because there are so many different approaches to it, start with Chap. 2. Both groups only need to skim the definitions—just come back later whenever you need them (Chap. 3). If you are an intermediate network analyst, i.e., you have conducted at least 3 network analytic projects, start with Chap. 8 and then read Chaps. 5–7. When you are an expert, just read the literacy chapters, starting from Chap. 10.

You will find that in this book, I often switch between the male and the female pronoun as long as I refer to some group of people in general (“the user → she” or “the user → he”). You find this annoying? Well, so do I! But as long as you and I notice it and still find it surprising or annoying or pleasing or anything else but normal, I feel the need to stick to it. Of course, the pronoun ‘she’ refers to both, male and female, persons. ☺

And thanks go to...

I would like to thank my colleagues Ulrik Brandes, Johannes Glückler, Kai Fischbach, and Alexander Mehler for our long discussions on network analysis. I would also like to thank the countless reviewers and foremost my own students Emöke-Ágnes Horvát, Wolfgang Schlauch, Mohammed Abufouda, and Sude Tavassoli for their influence on my work and the successful collaboration; last but not least, I would like to thank my collaborators from biology, especially Kevin Bähner and Thorsten Stoeck.

I hope that the book will foster a discussion on a more principled way of when to use which network analytic methods. The set of guidelines enabling this choice is what allows *network analysis*. However, this book is only the beginning of this and far from complete. Let me know your opinion, send in good and bad examples of network analysis, propose your own set of guidelines, and share all of this with me at zweig@cs.uni-kl.de. I will discuss a selection of those proposals on my blog <http://netz-werker.blogspot.de/>.

The book is dedicated to my husband who has shared all of my ups and downs in network analysis and supported me to write this book. Thanks for all the discussions on this topic that others not involved in network analysis might not have found as worthwhile as you.

Kaiserslautern
June 2016

Katharina A. Zweig

Contents

Part I Introduction

| | | |
|----------|---|----|
| 1 | A First Encounter | 3 |
| 1.1 | Introduction to Network Analysis | 3 |
| 1.2 | Data | 5 |
| 1.2.1 | From Relationship to Graph | 6 |
| 1.2.2 | First Probes into the Data | 8 |
| 1.2.3 | Measuring Indirect Effects | 12 |
| 1.2.4 | Distributions | 13 |
| 1.3 | Network Analysis Literacy: A Primer | 15 |
| 1.3.1 | Visualizations | 15 |
| 1.4 | Approaches to Network Analysis | 18 |
| 1.5 | Outlook | 19 |
| 1.6 | Recommended Reading | 20 |
| | References | 21 |
| 2 | Graph Theory, Social Network Analysis, and Network Science | 23 |
| 2.1 | Introduction | 23 |
| 2.2 | The Basis | 24 |
| 2.2.1 | Graph Theory | 24 |
| 2.2.2 | The Origins of Social Network Analysis in Sociology | 27 |
| 2.2.3 | Typical Viewpoints of Social Network Analysis | 30 |
| 2.2.4 | Network Science | 31 |
| 2.3 | Universal Structures versus Individual Features | 35 |
| 2.3.1 | Statistical Physics and Early Complex Network Analysis | 37 |
| 2.3.2 | Statistical Physics and Complex Network Analysis | 38 |
| 2.3.3 | Complex Network Analysis in Other Disciplines | 40 |

- 2.4 Network Analysis Literacy: General Requirements 42
 - 2.4.1 Implementations and Verbal Descriptions of Network Analytic Measures: A Primer 42
 - 2.4.2 Interpreting a Measure’s Value: A Primer 43
 - 2.4.3 Interpretation by Trained Domain Experts 45
 - 2.4.4 Interpretation by Academic Experts 48
 - 2.4.5 The Widespread Use of Scientific Rituals 49
 - 2.4.6 The Interpretation of Network Analytic Measures 49
- 2.5 Recommended Reading 53
- 2.6 Exercise 53
- References. 53
- 3 Definitions 57**
 - 3.1 Introduction 57
 - 3.2 Mathematical Abbreviations 58
 - 3.3 Set Theoretic Terms 58
 - 3.3.1 Function 60
 - 3.3.2 Partitions and Hierarchical Clustering 60
 - 3.4 Mathematical Operators 61
 - 3.5 Graph Theoretic Definitions 61
 - 3.5.1 Distances in Graphs 63
 - 3.5.2 Degrees and Walks in Graphs 63
 - 3.5.3 Graph Families 65
 - 3.6 Data Structures for Graphs 66
 - 3.6.1 Basic Data Structures. 67
 - 3.6.2 Basic Data Structures for Simple Graphs. 68
 - 3.6.3 Data Structures and Definitions for Directed Graphs 71
 - 3.6.4 Weighted Graphs 72
 - 3.6.5 Bipartite and Affiliation Networks 73
 - 3.6.6 Multiplex Networks 74
 - 3.7 Graph File Formats 74
 - 3.7.1 Graph Formats for Visualization 77
 - 3.8 A Little Bit of Linear Algebra 77
 - 3.8.1 Scalar Product 77
 - 3.9 Normalization 78
 - 3.9.1 Covariance. 78
 - 3.9.2 Correlation Coefficient. 79
 - 3.10 Algorithms and Runtime Complexity 80
 - 3.11 Plots and Diagrams. 81
 - 3.12 Distributions 82
 - 3.13 A Bit of Statistics 82
 - 3.14 Markov Chains 83
 - 3.14.1 Properties of Markov Chains 85

| | | |
|------|---------------------------|----|
| 3.15 | Further Reading | 86 |
| 3.16 | Exercises. | 86 |
| | References. | 88 |

Part II Methods

| | | |
|----------|---|------------|
| 4 | Classic Network Analytic Measures | 91 |
| 4.1 | Introduction | 91 |
| 4.2 | Direct Statistics. | 92 |
| 4.3 | Distance Based Measures | 93 |
| 4.4 | Degree Based Measures | 95 |
| | 4.4.1 Degree Distribution | 95 |
| | 4.4.2 Assortativity | 95 |
| 4.5 | Mutuality, Transitivity, and the Clustering Coefficient | 99 |
| | 4.5.1 Mutuality or Reciprocity | 99 |
| | 4.5.2 Transitivity | 100 |
| 4.6 | Density | 102 |
| 4.7 | Summary | 104 |
| 4.8 | Further Reading | 105 |
| 4.9 | Exercises. | 105 |
| | References. | 107 |
| 5 | Network Representations of Complex Systems | 109 |
| 5.1 | Introduction | 109 |
| 5.2 | Why Networks are only Models of Complex Systems | 109 |
| | 5.2.1 Edges as Abstract Representations of Real-World Relationships | 111 |
| | 5.2.2 Types of Network Representations | 113 |
| 5.3 | Phases of a Network Analytic Project. | 117 |
| | 5.3.1 Trilemma of Complex Network Analysis | 119 |
| 5.4 | Defining the Entity of Interest. | 121 |
| | 5.4.1 Network Boundary | 122 |
| | 5.4.2 Observing Entities | 123 |
| | 5.4.3 Entity Resolution. | 126 |
| 5.5 | Relationships and Mathematical Relations | 127 |
| | 5.5.1 Classic Relationships Analyzed in Complex Networks | 130 |
| 5.6 | Weighted and Dynamic Graphs | 131 |
| | 5.6.1 Observing and Representing Weighted Relationships | 131 |
| | 5.6.2 Dynamic Networks | 132 |
| | 5.6.3 Transformation into Undirected, Unweighted Networks | 133 |
| 5.7 | One-Mode Projections of Bipartite Graphs | 137 |
| | 5.7.1 Classic One-Mode Projections. | 137 |

- 5.7.2 Show Case: Co-authorship Networks. 139
- 5.8 An Example: Metabolic Networks 141
- 5.9 Summary 145
- 5.10 Further Reading 145
- 5.11 Exercises. 146
- References. 147
- 6 Random Graphs and Network Models 149**
 - 6.1 Introduction 149
 - 6.2 The Set of All Graphs with the Same Number of Nodes 150
 - 6.2.1 The $G(n,m)$ Random Graph Model 152
 - 6.3 The Classic Random Graph Model. 154
 - 6.4 The Small-World Model: Explaining the Small-World Phenomenon 158
 - 6.4.1 The Small-World Model (WS-Model) 162
 - 6.5 The Preferential Attachment Model (BA-Model) 165
 - 6.5.1 Scale-Freeness 166
 - 6.6 When is a Random Graph Model Explanatory? 170
 - 6.7 Summary 174
 - 6.8 Further Reading 175
 - 6.9 Exercises. 177
 - References. 179
- 7 Random Graphs as Null Models 183**
 - 7.1 Introduction 183
 - 7.2 Assessing the Significance of a Structural Feature 183
 - 7.2.1 Reciprocity Revisited I 184
 - 7.2.2 What is the Best Null Model for Assessing Reciprocity in General? 186
 - 7.2.3 Node Similarity and Co-occurrence. 187
 - 7.3 Fixed and Expected Degree Sequence Models 191
 - 7.3.1 Stub or Configuration Method. 193
 - 7.3.2 Simple Independence Model (SIM)—Approximating the Configuration Model 194
 - 7.3.3 Chung-Lu-Model: Expected Degree Sequences 196
 - 7.3.4 Fixed Degree Sequence Model 197
 - 7.4 The Philosophy behind Identifying Statistically Significant Structural Features. 198
 - 7.5 History of Assessing the Significance of Real-World Network Structures 201
 - 7.5.1 Network Motifs 202
 - 7.5.2 The Algorithm. 203
 - 7.5.3 Biologically Meaningful Motifs. 206
 - 7.5.4 Choosing the Best Null Model 207

| | | |
|----------|---|------------|
| 7.6 | Summary | 208 |
| 7.7 | Further Reading | 209 |
| 7.8 | Exercises. | 210 |
| | References. | 213 |
| 8 | Understanding and Designing Network Measures. | 215 |
| 8.1 | Introduction | 215 |
| 8.2 | Beware of verbal Descriptions—Why Mathematical Equations are Necessary | 216 |
| 8.2.1 | Reciprocity | 218 |
| 8.3 | Profile of a Measure’s Behavior | 221 |
| 8.3.1 | Applicability | 222 |
| 8.3.2 | Range of the Measure and Extremal Graphs | 225 |
| 8.3.3 | Scalability | 226 |
| 8.3.4 | Size Independence/Comparability | 227 |
| 8.3.5 | Robustness. | 228 |
| 8.3.6 | Assumptions | 228 |
| 8.4 | How to Design a Network Analytic Measure | 230 |
| 8.4.1 | Generalizing a Method | 231 |
| 8.4.2 | Another Interpretation of the Degree in Weighted Graphs. | 235 |
| 8.4.3 | Clustering Coefficient for Bipartite Graphs | 235 |
| 8.5 | Summary | 238 |
| 8.6 | Recommended Reading | 238 |
| 8.7 | Exercises. | 239 |
| | References. | 241 |
| 9 | Centrality Indices | 243 |
| 9.1 | Introduction | 243 |
| 9.2 | What is a Centrality Index? | 244 |
| 9.3 | Classic Centrality Indices | 246 |
| 9.3.1 | Degree-Like Centralities | 246 |
| 9.3.2 | Closeness-Like Centralities | 249 |
| 9.3.3 | Stress and betweenness-Like Centralities. | 250 |
| 9.3.4 | Correlation between Different Centrality Indices | 255 |
| 9.3.5 | Comparing Centrality Values in Different Networks. | 256 |
| 9.3.6 | The Centralization of a Graph | 258 |
| 9.4 | Generalizing Centrality Indices. | 259 |
| 9.4.1 | Centrality Indices for Networks between Different Groups of Nodes. | 259 |
| 9.4.2 | Centrality Indices for Directed Networks. | 260 |
| 9.4.3 | Centrality Indices for Weighted Networks. | 260 |
| 9.5 | Characterizations of Centrality Indices | 261 |
| 9.5.1 | The Graph-Theoretic Perspective. | 261 |
| 9.5.2 | Network Flow Processes and Centrality Indices | 264 |

| | | |
|-------|---|-----|
| 9.6 | Centrality-Based Visualization of Graphs | 264 |
| 9.7 | Applications of Centrality Indices | 265 |
| 9.7.1 | Centrality Distributions as General Structural Descriptors | 266 |
| 9.7.2 | Correlation between Centrality Indices and External Properties | 268 |
| 9.7.3 | Centrality Indices as Process-Based Predictors | 270 |
| 9.8 | Summary | 271 |
| 9.9 | Further Reading | 271 |
| 9.10 | Exercises | 272 |
| | References | 274 |

Part III Literacy

| | | |
|-----------|--|------------|
| 10 | Literacy: Data Quality, Entities, and Nodes | 279 |
| 10.1 | Introduction | 279 |
| 10.2 | Describing a Network Representation Transparently | 280 |
| 10.3 | Bad Data | 282 |
| 10.3.1 | Bad Data: Protein-Protein Interaction Networks | 282 |
| 10.3.2 | Bad Data: BGP Routing Data | 286 |
| 10.3.3 | Inferred Transcription Network Data | 287 |
| 10.4 | Network Boundary | 289 |
| 10.4.1 | When is a Node a Node | 289 |
| 10.5 | Sampling Effects | 293 |
| 10.5.1 | Dynamic and Time-Thresholded Data | 296 |
| 10.6 | Evaluating Sampling Strategies | 297 |
| 10.6.1 | Evaluating BGP/Traceroute Data | 298 |
| 10.7 | Data Biases | 299 |
| 10.7.1 | Data Biases in Protein-Protein Interaction Data | 299 |
| 10.7.2 | Data Biases in Surveys | 300 |
| 10.7.3 | Estimating the Degree of a Node in a Network | 302 |
| 10.8 | Curating Complex Networks | 304 |
| 10.9 | Summary | 306 |
| 10.10 | Further Reading | 306 |
| 10.11 | Exercises | 306 |
| | References | 309 |
| 11 | Literacy: Relationships and Relations | 313 |
| 11.1 | Introduction | 313 |
| 11.2 | When is an Edge an Edge? | 314 |
| 11.3 | Aggregations in Time and Space | 318 |
| 11.3.1 | Aggregation in Time | 318 |
| 11.3.2 | Aggregation in Space | 319 |
| 11.3.3 | Choosing an Appropriate Observation Period | 320 |

- 11.4 Weighted Relationships. 322
 - 11.4.1 Interrelationship with Chosen Method 322
 - 11.4.2 Dynamic Weights 325
 - 11.4.3 Thresholding 326
- 11.5 Proxy Relationships 327
 - 11.5.1 Proxies for Sexual Relationship Networks. 327
 - 11.5.2 Online Social Network Data as Proxies. 329
 - 11.5.3 With Whom do We Discuss Important Matter. 330
 - 11.5.4 Co-authorship versus Collaboration. 332
 - 11.5.5 Interchangeability of Social Relations 332
 - 11.5.6 Observational versus Recalled Interactions 334
 - 11.5.7 Email Interaction versus Communication
Networks 334
 - 11.5.8 Internet Network Data and Their Proxies. 336
- 11.6 Relations that don't Lend Themselves to a Network
Representation 338
 - 11.6.1 Information Contained in Relations. 338
 - 11.6.2 Mathematical Relations without Network
Processes 340
 - 11.6.3 Aggregating Paths into Complex Networks. 341
 - 11.6.4 Relationships, Network Processes, and Complex
Networks 344
- 11.7 Horizons of Network Processes 348
- 11.8 Data Responsibility. 350
 - 11.8.1 Evaluating Existing Network Data for Re-use. 351
 - 11.8.2 Data Hygiene, Producer and Consumer Rules. 353
 - 11.8.3 Producer Rules: Making Data Reusable. 354
 - 11.8.4 Consumer Rules: Validating Data 356
- 11.9 Aim of Analysis (A-Rules) 357
 - 11.9.1 Publishers' Responsibility 357
- 11.10 Summary 358
- 11.11 Further Reading 359
- References. 359
- 12 Literacy: When Is a Network Model Explanatory? 363**
 - 12.1 Introduction 363
 - 12.2 Models of Networks and Processes. 365
 - 12.2.1 What is a Scientific Model?. 366
 - 12.2.2 Modelling Processes on Complex Networks 370
 - 12.2.3 Evolution of Models 371
 - 12.3 Structure, Function, and Behavior of Network Models. 372
 - 12.3.1 Interpretation of 'Smallness' as a Function 373
 - 12.3.2 Properties and Behavior of "Scale-Free"
Networks 377

- 12.4 Explanatory Models 381
 - 12.4.1 When Preferential Attachment is not Enough 382
 - 12.4.2 Networks with a “Scale-Free” Degree Distribution
Which are not “Scale-Free” 383
 - 12.4.3 The Internet—A “Scale-Free” Network without a
Hub-Dominated Architecture 384
 - 12.4.4 Shrinking Diameters in the Evolution of Complex
Networks 385
 - 12.4.5 Measuring Preferential Attachment 385
- 12.5 Summary 387
- 12.6 Further Reading 389
- References. 392
- 13 Literacy: Choosing the Best Null Model 395**
 - 13.1 Introduction 395
 - 13.2 Assessing the Small-World Phenomenon 398
 - 13.2.1 Clustering Coefficient in One-Mode Projections
of Bipartite Graphs 399
 - 13.3 The Rich-Club Coefficient 401
 - 13.4 Reciprocity Revisited II 405
 - 13.5 A New Perspective on One-Mode Projections 407
 - 13.5.1 The Simple Independence Model SIM. 408
 - 13.5.2 An Example: MovieLens. 410
 - 13.5.3 Discussion of the SIM. 415
 - 13.5.4 The Fixed Degree Sequence Model FDSM for
Bipartite Graphs. 418
 - 13.6 Evaluating Expectation Models by a Gold Standard or
Ground Truth 419
 - 13.6.1 Building the OMP. 420
 - 13.6.2 Is There a Weighted FDSM?. 421
 - 13.7 Can the Configuration Model Replace the FDSM?. 422
 - 13.8 Summary 425
 - 13.9 Further Reading 426
 - 13.10 Exercises. 427
 - References. 428
- 14 Literacy Interpretation 431**
 - 14.1 Introduction 431
 - 14.2 The Interpretation of Measures in the Context
of a Complex System 432
 - 14.3 Interpretation of Distance-Based Measures 435
 - 14.3.1 Robustness Measures Based on Distance. 435
 - 14.3.2 Comparing Average Distances of Different
Networks 440
 - 14.3.3 Interpretation of Low Average Distances
in Metabolic Networks 441

| | | |
|-----------|--|------------|
| 14.4 | Centrality Index Literacy | 443 |
| 14.4.1 | Borgatti’s Flow Concept | 444 |
| 14.4.2 | Interpretation of Classic Centrality Indices | 445 |
| 14.4.3 | Air Transportation Networks | 447 |
| 14.4.4 | Multiplex Air-Transportation Networks | 450 |
| 14.4.5 | Designing Interpretable Centrality Indices | 454 |
| 14.5 | Explorative Applications of Distance Based Measures | 455 |
| 14.6 | The Centrality of Moscow in the 12th and 13th Century | 457 |
| 14.7 | Sexual Contact Networks | 462 |
| 14.7.1 | From Data to Network. | 463 |
| 14.7.2 | The Human Web of Sexual Contacts. | 463 |
| 14.7.3 | An Assessment of Preferential Attachment as a Mechanism for Human Sexual Network Formation | 466 |
| 14.8 | Post-Hoc Analysis. | 467 |
| 14.9 | Verbal Description of Findings. | 469 |
| 14.10 | Summary | 470 |
| 14.11 | Exercises. | 471 |
| | References. | 472 |
| 15 | Ethics in Network Analysis. | 475 |
| 15.1 | Why Ethical Network Analysis Needs Network Analysis Literacy. | 475 |
| 15.2 | The Wegman Report. | 476 |
| 15.2.1 | Discrediting a Scientist by Co-authorship-Network Analysis. | 476 |
| 15.3 | Who Owns a Relationship?. | 479 |
| 15.4 | Prediction Based on Network Analysis. | 482 |
| 15.5 | Summary | 483 |
| | References. | 484 |
| | Appendix A: The Structure and Typical Outlets of Network Analytic Papers. | 487 |
| | Appendix B: Glossary. | 493 |
| | Appendix C: Solutions | 499 |
| | Author Index. | 529 |
| | Subject Index. | 531 |

Part I

Introduction

What is network analysis about?

Chapter 1

A First Encounter

Abstract The first chapter of the book gives a short overview of what network analysis does and why it is considered to be a vital part of complex system science: the network analytic framework allows to represent the interaction structure of a complex system as a complex network. The network's structure can then be analyzed by the application of several structural measures. However, there are two different branches in network analysis that either use the resulting values to find so-called *universal features* of complex systems or to allow a *contextual, semantic analysis*. The latter focuses on the connection between structure and function of a network with respect to the complex system of interest and some specific research question. There is a caveat, though: while, in principle, structural measures can be applied to all kinds of networks, if one is only searching for universal features, their results are not always interpretable with respect to a predefined research question. The term "network analysis literacy" is introduced to describe the knowledge of when to apply which measure to yield an interpretable result with respect to the complex system of interest.

1.1 Introduction to Network Analysis

Networks impress by their visual and intuitive quality: everyone of us is entangled in various friendship networks and business relationships, and the prospect of understanding the seemingly complex and erratic net of our personal relationships is an exciting one. Similarly, looking at scientific data in a new way, finding simple patterns that chip away individual noise to extract the main functional groups of entities in the complex system at hand, is surely one of the most gratifying moments in every scientist's life. Network analysis seems to be one of the most promising frameworks within which these two aspects, our personal life and our academic interest, can be combined, analyzed, and maybe even be understood. This prospect and the many exciting articles in journals such as *Science*, *Nature*, and *PNAS*, together with the interdisciplinary applicability of network analysis to various data sets and questions, has led to a tremendous interest in the methods provided by network analysis: Fig. 1.1 shows the dramatic increase of the number of articles with the keywords "network

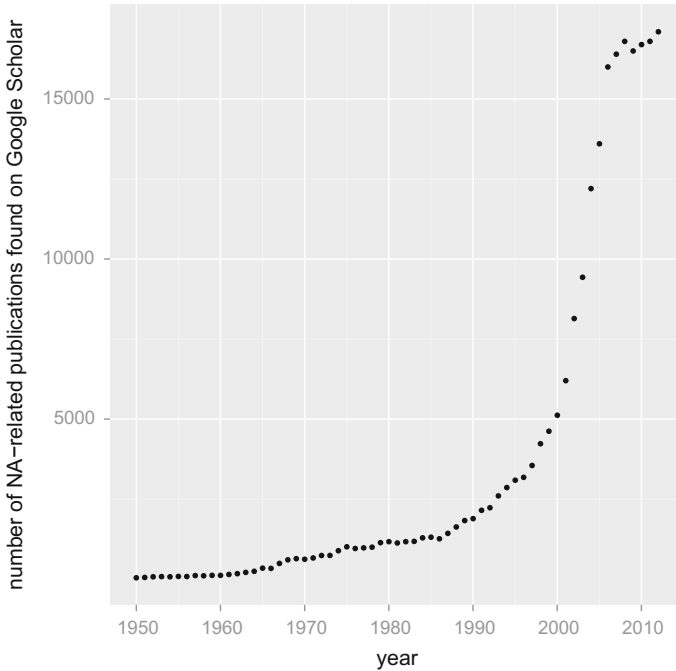


Fig. 1.1 Number of articles published in the given year containing the exact phrases “network analysis” and “complex network” as given by Google Scholar on the 12th of October, 2013

analysis” or “complex networks” as found by Google Scholar.¹ Starting from about 100 articles (as found on Google scholar) in the 1950s, the last years saw more than 15,000 articles with these terms.² This chapter gives a broad first encounter with network data by showing the first steps in analyzing a new set of network data.

The following chapters present a classic and widely used part of the toolkit in network analysis; but more importantly, they elaborate on the questions that are necessary to be answered in order to decide whether a given method is meaningful to the research question. One main caveat in network analysis is that once data is transformed into a graph representation, one can in principle apply any of the hundreds of network analytic methods to it—but not every method will compute meaningful and interpretable results with respect to the given data and the question

¹For each year from 1950 to 2012, a Google Scholar search with both terms, connected by an “OR” was conducted. The number of results displayed was taken as the data point for the given year. The number of results is unlikely to hit the number of published articles in any way but gives at least an indication of the strongly increased interest in the topic.

²Note that double counting is as likely as an underestimation of the number of articles: articles with this topic may, for example, have been overlooked because they were published in a non-public journal which Google Scholar might not have access to. Again, the number given by Google scholar is only an indication of how many articles really have been published.

to answer. This book is thus not so much about introducing measures, many more can be found in the books by Wasserman and Faust [24], Newman [18], Borgatti et al. [6], or the book edited by Brandes and Erlenbach [7]. This book rather focuses on this last part, which I call *network analysis literacy*: it aims to empower its readers to know when to use which method so that they can quickly delve into the exciting analysis of networks themselves. Be warned that the technique by which this is accomplished follows the Socratic method which, in general, poses more questions than gives definite answers.

So, what is the first step in network analysis? One basic phase is the transformation of relational data into a complex network representation as described in the next section.

1.2 Data

The first question you might have is: what kind of data can actually be meaningfully represented as networks? A first answer is: almost any kind of data. The basic requirement is that there is a **distinct set of entities**, e.g., humans, organizations, proteins, computers, or books. The second requirement is that there is a known **relationship** between these entities. The information of whether any two entities are in the given relationship or not needs to be known for a large part of the entities since otherwise any kind of analysis will be quite shaky. Some obvious relationships between persons are: friendship, kinship, or employee-employer-relationships. Another interesting type of relationship is membership: it is a relationship between two different kinds of entities, namely persons and institutions, but that can also be represented by a network.

Relationships between non-human entities are equally abundant, and in many cases the resulting structures are also termed *networks* in our day-to-day language: examples are metabolic networks, protein-protein-interaction networks, neural networks, street networks, or computer networks. All of the above examples might be considered ‘natural networks’ but are there more abstract relationships that can also be represented as complex networks?

Interestingly, mathematicians have a very general understanding about what is a relation and what is not: in a mathematical sense two books can be defined to be “related” because their cover was created by the same designer. This “relatedness” does not mean that they are necessarily related in any colloquial sense: their content can of course be very different! Nonetheless, in a mathematical sense, the relation is meaningfully defined and can be easily checked by an external observer. In mathematics, a *relation* is simply defined as a subset of pairs of entities: $R \subseteq O \times O$, where $O \times O$ denotes the set of all possible pairs from some set of entities or objects O .

Note 2. Mathematically, a *relation* R on a given set of entities or objects is just an arbitrary choice of pairs of these entities (objects), denoted by $R \subseteq O \times O$. In principle, any relation can be represented as a graph.

This is on the one hand much more general than the day-to-day notion of a relationship but on the other hand much less intuitive: a *relation* does not need to stand in any correlation to a real-world *relationship*; it can even represent a relationship that would not be seen as meaningful in the real world: for example, all humans with the same first name can be represented by a relation or all humans which share the same last digit of their ID-card number. Mathematical relations can also (meaningfully) be derived from other relations: One can build a second network based on the connection structure of another network by, for example, connecting two persons with each other if they share at least 8 friends in a friendship network. In this second network, there might be two persons that are connected because they share enough friends but which are not befriended themselves, and vice versa.

So, in this book, a *relationship* is something that can be observed in the real world, a *relation* is the mathematical structure which possibly represents a relationship. However, not all relations are associated with any relationship and the same relation, i.e., the same subset of pairs of a given set of elements, can represent different relationships. How is now a relationship turned into a complex network? This is discussed in the following.

1.2.1 From Relationship to Graph

In the moment a set of entities and a relationship of interest has been defined, there is a range of decisions to be made, to turn the concept of that relationship into a procedure that decides, for each pair of entities, whether they are in the associated mathematical relation or not. In most cases, when data is turned into a network representation, several decisions have to be made: if the relationship of interest contains a direction, is it necessary to include this information in the mathematical relation? Are there different levels of intensity of a given relationship and is it necessary to differentiate between them, by assigning weights to the pairs in the relation? Each of these decisions changes the set of available structural measures and the interpretation of the measure if it is applied to the network. Chapter 5 will explain in detail how data can be turned into networks. In any case, mathematically, a graph is the combination of a set of elements and a relation defined on these elements.

Note 3. What is the difference between a (complex) *network* and a *graph*? The quick answer is that a graph is the *abstract representation* of a relation between entities while a network combines the graph with additional information about the entities and their **relationship** represented by the graph.

In most cases, a *complex network* represents only one set of entities (sometimes two) and **one relationship** between the entities, with some limited options on the attributes assigned to the (mathematical) relation and usually no attributes assigned to the elements. On the graph level, the elements are called *nodes* or *vertices*, and the pairs of nodes contained in the relation are called *edges*. The graph can indicate whether the relationship is directed by either containing a symmetric or asymmetric relation. In the first case, whenever (a, b) is contained in the relation, so is (b, a) : for example, the graph can store information of whether Tim is father of Tom or vice versa (or none of that). It can also represent weights of the relationship, by assigning a weight to each element in the relation. Again, the graph itself does not store information about which entity is represented by which node, it is oblivious of any identity of the node. Thus, the *graph* is the more abstract representation which mainly concentrates on the connection structure.

The *network* makes the connection between the graph and the complex system whose interactions it represents. Especially, the network assigns entities to nodes. Furthermore, it is the set of all descriptions and observations of the system in which the entities and their relationship is valid. It can contain observations on the entities like the age and gender of a human actor or the year of publication of a film. It can also contain more than one type of relationship between the actors, or additional observations about the relationships between entities, like the duration of all calls between mobile phone users.

In summary, a *complex network* is a graph, in which the set of elements is associated with a set of entities, and in which the relation between these elements represents a relationship between the corresponding entities.

The distinction between a network and its graph is often not very important, and thus *network* and *graph* are used quite interchangeably in most articles and also in this book.

Note 4. The promise of network analysis is that the abstraction of a complex system as represented by a complex network and its underlying graph still allows to infer something about the complex system of interest. That is actually a strong assumption and later chapters (e.g., Chaps. 10 and 14) elaborate preconditions to enable this transfer.

So, what are the first steps after the data is represented as a network? The following section shows some typical approaches to get a first impression of a new data set on the example of a movie-co-rating network.

Table 1.1 The movie-movie-similarity network contains 494 films and represents 9796 relationships between them

| Statistics | Value |
|------------|-------|
| n | 494 |
| m | 9796 |
| $\rho(G)$ | 0.08 |

1.2.2 First Probes into the Data

The movie-co-rating network is a network deduced from a so-called *bipartite graph*: As indicated above, some data sets describe a relationship between two different entities, for example, how customer of a video rental store rate the films they rented. Such a data set documents a relationship between customers and the films they rated, but there is no direct relationship between any two customers or between any two films. This kind of data is represented by a bipartite graph, that is, one that can be split in two parts such that all known relationships are only between entities from different parts. Based on such data, one can compute a similarity measure between the films that quantifies whether the films have been more often liked by the same persons than expected or not; this technique is called a *one-mode projection* of a bipartite graph and described in Sect. 13.5.

Such a one-mode projection is the basis for the following demonstrations. It has been created such that the relation is undirected and it is assumed that a pair of films connected by an edge are similar in content. The data comes in a format³ that is readable by various software applications, e.g., Gephi⁴ which is well suited for visually exploring a medium-sized graph [22]. Similarly well suited are yEd [27], visone [13], or Cytoscape [26].

The very first useful information about the data is how many films it contains and how many relationships between them exists. In most visualizations of graphs, this information is given immediately when the graph is opened and displayed. In general, the number of entities is denoted by n and the number of relationships is denoted by m . With around 500 nodes and 10,000 edges, the network is of medium size. From these two basic statistics the so-called *density* of relationships can be computed as another, first inspection into the graph. It is defined as the number of existing relationships divided by the number of possible relationships: in principle, every pair of entities could be related to each other, thus, the number of possible relationships is given by $n(n - 1)/2$. The density in the given data with $n = 494$ and $m = 9796$ can thus be computed to be 0.08. This density can also be interpreted as the *probability* that a randomly chosen pair of movies is related; this probability is obviously very small. Table 1.1 summarizes the basic statistics.

³Sections 3.6 and 3.7 discusses various graph data formats and how they can be transformed into each other.

⁴Freely downloadable from <http://gephi.org/>.